
COMPARATIVE ANALYSIS AND EVALUATION OF AGING FORECASTING METHODS FOR SEMICONDUCTOR DEVICES IN ONLINE HEALTH MONITORING

Adrian Villalobos

Electronics & Computer Science Department
Mondragon University
Mondragon, 20500 (Spain)
adrian.villalobos@alumni.mondragon.edu

Iban Barrutia

Electronics & Computer Science Department
Mondragon University
Mondragon, 20500 (Spain)
ibarrutia@mondragon.edu

Rafael Peña-Alzola

Electronics & Electrical Engineering Department
University of Strathclyde
Glasgow, G11RD (UK)
rafael.pena-alzola@strath.ac.uk

Tomislav Dragicevic

Wind & Energy Systems Department
Technical University of Denmark (DTU)
Lyngby, 2800 (Denmark)
tomdr@dtu.dk

Jose I. Aizpurua

Computer Science and Artificial Intelligence Department, University of the Basque Country (UPV/EHU)
Ikerbasque, Basque Foundation for Science
San Sebastian, 20018 (Spain)
joxe.aizpurua@ehu.eus

March 27, 2025

ABSTRACT

Semiconductor devices, especially MOSFETs (Metal-oxide-semiconductor field-effect transistor), are crucial in power electronics, but their reliability is affected by aging processes influenced by cycling and temperature. The primary aging mechanism in discrete semiconductors and power modules is the bond wire lift-off, caused by crack growth due to thermal fatigue. The process is empirically characterized by exponential growth and an abrupt end of life, making long-term aging forecasts challenging. This research presents a comprehensive comparative assessment of different forecasting methods for MOSFET failure forecasting applications. Classical tracking, statistical forecasting and Neural Network (NN) based forecasting models are implemented along with novel Temporal Fusion Transformers (TFTs). A comprehensive comparison is performed assessing their MOSFET ageing forecasting ability for different forecasting horizons. For short-term predictions, all algorithms result in acceptable results, with the best results produced by classical NN forecasting models at the expense of higher computations. For long-term forecasting, only the TFT is able to produce valid outcomes owing to the ability to integrate covariates from the expected future conditions. Additionally, TFT attention points identify key ageing turning points, which indicate new failure modes or accelerated ageing phases.

Keywords Semiconductor · Prognostics · Forecasting · Condition monitoring · Temporal Fusion Transformer · Neural Networks

1 Introduction

The increased deployment of renewable energy plants has led to an increased research activity focused on prognostic models for semiconductor-based power modules as a way to reduce the operation and maintenance costs through predictive maintenance (Zhu et al., 2024). The main goal of failure prognostics is the reliable and accurate remaining useful life (RUL) forecast of the component under study by modelling the analysed failure mode and forecasting the future ageing evolution (de Pater and Mitici, 2023).

Semiconductors are used as electronic switches due to their high efficiency and fast switching properties (Zhang et al., 2023). However, they are often ranked as having the lowest reliability in power systems (Novak et al., 2021). The main failure modes of semiconductors can be classified into sudden and ageing failures (Hanif et al., 2019). Sudden failures are caused by random phenomena such as cosmic radiation or electric discharge, e.g. (Kang et al., 2025). In contrast, aging failures are caused by environmental or operational stress that exceeds a failure threshold limit, e.g. (Zubizarreta et al., 2025). With regard to aging failures, the main focus of this research is on fatigue damage, which is caused by an imbalanced coefficient of thermal expansion of the different materials within the semiconductor structure.

MOSFETs are semiconductor devices used in different power applications (Lutz et al., 2018). The reliability and lifetime assessment of MOSFETs is complex and influenced by different aging processes that vary with cycling and temperature (Anderson et al., 2025). The main failure mechanism is die-attachment degradation, which leads to increased thermal impedance and higher device temperature, finally deriving in bond wire lift-off (Celaya et al., 2012). Thus, the degradation of lead-free solder die-attach results in an increase of the activation resistance or on-state, $R_{DS_{ON}}$, which depends on the junction temperature. Therefore, monitoring the evolution of $R_{DS_{ON}}$ can be used to determine the RUL of the device under thermal stress. Namely, the MOSFET die attachment degradation cycle occurs due to the repeated operation of the MOSFET, which transits between the ON and OFF states and this causes an increase of the thermal impedance and device temperature, and eventually the occurrence of bond wire lift-off. Figure 1 shows an example of the MOSFET degradation cycle, estimated from the dataset used in the case study (cf. Section 4). The drain current I_D is used to estimate the $R_{DS_{ON}}$, which is the main precursor to predict the bond wire lift-off.

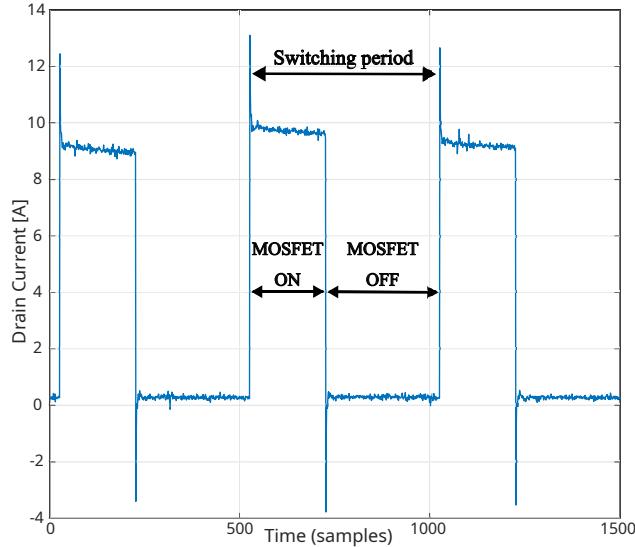


Figure 1: MOSFET degradation cycle example calculated from the case study dataset (cf. Section 4).

1.1 Related Works

Artificial intelligence (AI) and machine learning (ML) applications are being widely developed for health monitoring applications for power electronic devices, e.g. design for reliability of power electronic systems (Dragičević et al., 2019) or parameter design for converters (Wang et al., 2024). A comprehensive review of AI applications in power electronics is provided by (Zhao et al., 2021). Among these applications, MOSFET prognostics have received significant attention, with NASA's MOSFET degradation dataset serving as a common benchmark (Celaya et al., 2012) – see Section 4 for more details. Prognostics for MOSFETs rely on key failure precursors including $R_{DS_{ON}}$, drain to

source voltage V_{DS} , and junction temperature T_j . Various methodologies, including statistical forecasting models, state-space approaches, and ML methods have been explored to enhance the accuracy of RUL prediction.

Early studies used time-series and regression-based approaches to forecast degradation trends. (Alonso et al., 2019) analysed the deviation of $R_{DS_{ON}}$ from its pristine condition as a health indicator for RUL estimation. Under the assumption of an exponential degradation model, they explored Least Squares, Horizontal Average, and their linearized versions to improve one-step-ahead predictions. Similarly, (Wang et al., 2022) compared ARIMA and Gauss-Newton nonlinear regression, evaluating their accuracy in predicting $\Delta R_{DS_{ON}}$. Their results demonstrated that ARIMA provided more precise RUL estimates at different degradation stages.

With the rise of data-driven methods, different ML models have demonstrated improvements in predictive accuracy. (Li et al., 2018) integrated Echo State Networks (ESN) with Particle Filters (PF) to model the degradation process. The ESN parameters were dynamically updated using $R_{DS_{ON}}$ as a precursor. In an alternative approach, (Baharani et al., 2019) proposed a cloud-based monitoring framework using Long Short-Term Memory (LSTM) networks to track the variation of $R_{DS_{ON}}$, showing higher accuracy than Kalman and Particle Filters. Building on this, (Kang et al., 2024) developed an LSTM-based prognostic approach, trained and validated on NASA's dataset with different training/testing proportions to optimize regression performance.

Beyond deep learning, alternative ML approaches have emerged for enhanced feature extraction and predictive reliability. (Kahraman et al., 2024) designed an ML pipeline integrating feature extraction, classification, and Bayesian Ridge Regression (BRR) for RUL prediction. RUL estimation was performed only if the classifier diagnosed the MOSFET in a pre-failure state. Feature-based methodologies have also proven effective in identifying MOSFET degradation patterns. (Yang et al., 2022) compared Nonlinear Auto-Regressive Model with Exogenous Inputs (NARX) and Partial Least Squares (PLS) regression, both combined with a Cumulative Sum (CUSUM) analysis. Their approach demonstrated superior accuracy compared to Principal Component Analysis (PCA) and exponential weighted moving average (EWMA-PCA) in detecting early-stage degradation.

Stochastic models offer another perspective on MOSFET prognostics by incorporating uncertainty in failure evolution. (Zhao et al., 2018) formulated MOSFET degradation as a Continuous-Time Markov Chain (CTMC), integrating it into Cox's Proportional Hazards (PH) model to estimate time-to-failure. Here, filtering and discretizing $R_{DS_{ON}}$ measurements using k -means clustering allowed for confidence-interval-based RUL predictions. Addressing the inherent variability in failure thresholds, (Wu et al., 2024) proposed a nonlinear Wiener process, continuously updating model parameters until reaching a probabilistic failure threshold, enabling a more adaptable prediction framework.

Beyond MOSFETs, similar prognostics methods have been applied to Insulated Gate Bipolar Transistors (IGBTs), using different health indicators as failure precursors: collector-emitter voltage, V_{ce} , collector current, I_c , junction temperature, T_j , or transient thermal impedance, TI . These precursors capture degradation mechanisms related to gate oxide wear out, bond wire degradation, and solder fatigue, which are common aging failure modes in IGBTs (Zubizarreta et al., 2025). Early approaches focused on time-series modeling and statistical estimation. (Alghassi et al., 2016) proposed an RUL estimation framework based on Temporal Delayed Neural Networks (TDNNs). The model took four sequential measurements of V_{ce} and their variations to estimate the component's health state. A moving average filter was applied to remove noise, and a Normal distribution assumption helped refine predictions. This health-state was then combined with Maximum Likelihood Estimation (MLE) to infer confidence intervals for the RUL, offering a probabilistic failure prognostics.

Unsupervised learning techniques have also been explored to capture degradation patterns without labeled failure data. (Rigamonti et al., 2018) developed an ensemble of Self-Organizing Maps (SOMs) for component- and population-based degradation state identification. The ensemble dynamically weighted population-based SOMs providing general degradation trends with component-specific SOMs tailored to individual device behavior. Evaluations on an experimental dataset containing V_{ce} , I_c , and case temperature T_c measurements showed reduced classification errors compared to static aggregation and single-component SOM models.

Deep learning architectures have been also used to extract complex IGBT degradation features. (Xiao et al., 2022) introduced a self-attention-based Neural Network (NN) trained on a dataset containing T_j , I_c , V_{ce} , and package temperature. Instead of directly using raw sensor data, the study post-processed the signals to compute transient thermal impedance (TI), a direct failure precursor for IGBTs. 21 handcrafted statistical features were extracted and refined to form the dataset for improved prognostics. The model was evaluated under both offline and online training/testing strategies, demonstrating its adaptability for real-world applications.

State-of-the-art graph-based and physics-informed ML approaches have recently gained traction in IGBT prognostics. (Deng et al., 2024) introduced a Spatio-Temporal Fusion Graph Network (STFGN), where IGBT health states were modeled as graphs constructed from discrete voltage levels, enabling a structured approach to degradation prediction. Meanwhile, (Lu et al., 2023) applied Physics-Informed Neural Networks (PINNs) for RUL estimation, integrating

domain knowledge into the Neural Network's loss function. This approach enforced monotonic degradation constraints and set initial health state values, ensuring physically meaningful predictions. Expanding on this direction, (Fassi et al., 2024) provided an extensive review of power converter maintenance strategies based on physics-informed ML methods.

Table 1 displays the synthesis of related work organised according to semiconductor technology, ageing precursor, methods, prediction horizon, and whether the approach focuses on forecasting or regression models (F/R).

Table 1: Summary of recent semiconductor ageing forecasting works.

Reference	Technology	Precursor	Methods	Horizon	F/R
(Alonso et al., 2019)	MOSFET	$R_{DS_{ON}}$	LS variants	1 step ahead	F
(Li et al., 2018)	MOSFET	$R_{DS_{ON}}$	ESN-PF	70%-30%, LOO	F, R
(Wang et al., 2022)	MOSFET	$R_{DS_{ON}}$	ARIMA, Gauss-Newton	$t_p = [100, 120, 140, 160, 180, 200]$ min	F
(Yang et al., 2022)	MOSFET	$R_{DS_{ON}}, I_D, V_{DS}$	NARX/PLS-CUMSUM	Anomaly Det.	-
(Zhao et al., 2018)	MOSFET	$R_{DS_{ON}}$	CTMC-PHs	Different times [min]	R
(Kahraman et al., 2024)	MOSFET	$R_{DS_{ON}}$	BRR	Different instants from pre-failed state	R
(Wu et al., 2024)	MOSFET	$R_{DS_{ON}}$	Wiener Process	Continuous parameter update	R
(Kang et al., 2024)	MOSFET	$R_{DS_{ON}}$	LSTM	train/test proportions on test #36	R
(Baharani et al., 2019)	MOSFET	$R_{DS_{ON}}$	LSTM	104 samples	R
(Alghassi et al., 2016)	IGBT	V_{ce}	TDNN	Long-term	R
(Rigamonti et al., 2018)	IGBT	V_{ce}, I_{col}, T	SOMs	Classification	R
(Xiao et al., 2022)	IGBT	TI	Self-attention NN	Train/test 30, 50, 70%. LOO, 30%	F, R
(Deng et al., 2024)	IGBT	V_{ce}	STFGN	weak supervised training	R
(Lu et al., 2023)	IGBT	V_{ce}	PINN	Train/test 80%/20%, out-of-sample	R

Legend: F/R: forecast or regression; LOO: leave one out.

Table 1 displays that the main failure precursor of MOSFETs is the activation resistance, $R_{DS_{ON}}$. Regarding the used predictive methods, mainly statistical methods (Alonso et al., 2019; Wang et al., 2022; Yang et al., 2022; Zhao et al., 2018), ML methods (Baharani et al., 2019; Kang et al., 2024; Kahraman et al., 2024; Rigamonti et al., 2018; Alghassi et al., 2016; Xiao et al., 2022; Lu et al., 2023; Deng et al., 2024), and combination of state-space and ML methods have been employed (Li et al., 2018).

Most of the revised methods use regression models with covariates to inform the predictive model (Yang et al., 2022; Zhao et al., 2018; Alghassi et al., 2016; Baharani et al., 2019; Rigamonti et al., 2018; Kahraman et al., 2024), and only few of them focus on forecasting the future evolution of the ageing precursor using only past information. Regression models can predict the future ageing evolution, however, at the cost of monitoring additional variables and with limited ability to capture historical dependencies. In contrast, forecasting applications use only past information to predict the future aging evolution, which is more difficult, but also more cost-effective. Regression and forecasting methods mainly differ in time-dependency, incorporation of time-dependent properties such as autocorrelation or seasonality, and the treatment of covariates or exogenous variables.

Accordingly, the main focus of this work is to forecast the future evolution of the MOSFET aging trajectory based solely on past information. That is, this process requires learning historical degradation dynamics and forecasting future aging evolution, which is a challenging prediction task. Among the forecasting methods, most of the methods rely on one-step-ahead prediction models, which are sequentially integrated to predict longer horizons (Alonso et al., 2019; Li et al., 2018; Wang et al., 2022). Interestingly, (Xiao et al., 2022) presented an IGBT health index forecasting through feature extraction, model learning and RUL prediction. Different from regression models, forecasting models may present a cost-effective monitoring strategy, which may be on-boarded in real time applications due to the ability of making predictions based on a single signal.

1.2 Contribution, Objectives & Impact

Due to the rapid and widespread development of ML models for forecasting applications (Benidis et al., 2022; Lim and Zohren, 2021), the use of powerful new methodologies, such as transformer-based forecasting methods

(Vaswani et al., 2017), can generate useful insights for extended observations (Gilpin, 2023). They have been proven to be valid, especially for long-term predictions (Biggio et al., 2023). Unlike regression methods that consider multiple independent variables sampled at the same time instant of interest to explain the response variable, transformer-based forecasting models can incorporate covariates and historical dependencies.

Consequently, the main contribution of this work is the detailed comparative evaluation of classical forecasting methods with transformer-based forecasting methods, which have the ability to incorporate expected future covariates along with other influencing factors. To this end, first, it is necessary to adapt transformer models for MOSFET aging forecasting purposes through a fit-for-purpose architecture which uses forecasting covariates. Subsequently the capability of transformer models for short- and long-term forecasting has been exhaustively assessed and compared with a comprehensive benchmarking approach that includes different statistical, state-space based filtering, and ML methods.

There are works that analyze AI and ML methods for power electronic applications (Zhao et al., 2021). However, their focus is on the bibliometric research assessment and the identification of research gaps. In contrast, the goal of this research is the practical evaluation and critical comparison among failure forecasting methods for short- and long-term forecasting horizons through the implementation of classical statistical, state-space filtering and ML methods, along with state-of-the-art transformer-architecture based forecasting methods. The proposed comparative evaluation has a direct impact on the health management of semiconductors in general and MOSFETs in particular. On the one hand, the proposed evaluation framework identifies trade-offs among state-of-the-art forecasting methods for MOSFET failure forecasting. On the other hand, transformer-based architectures are adapted for MOSFET failure forecasting, which can greatly improve long-term MOSFET failure forecasting activities.

1.3 Organization

The remainder of this article is organized as follows. Section 2 reviews MOSFET ageing basics. Section 3 presents the proposed ageing forecasting assessment approach. Section 4 defines the case study. Section 5 applies the proposed approach to the case study. Section 6 provides a discussion of the presented methodology, and finally, Section 7 concludes.

2 MOSFET Degradation Modelling

MOSFETs work in the ohmic region when ON, and in the cut-off region when OFF. The current through the semiconductor, I_{DS} , subject to drain to the source voltage, V_{DS} , is controlled by a voltage applied between the gate and the source terminals. Even though the MOSFET gate is electrically isolated, the gate driver needs to supply sufficient current to charge all the capacitances across the device. In particular, the capacitance between the gate and the drain (also called reverse capacitance) increases the effective input capacitance seen at the gate by orders of magnitude due to the Miller effect. Once the MOSFET is fully turned on, the voltage drops across source and drain can be modelled as an on-resistance, $R_{DS_{ON}}$, whose value is dependent on the junction temperature.

The $R_{DS_{ON}}$ can be monitored by recording the values of the drain current I_D and the drain to source voltage V_{DS} in the device under ageing circumstances. On the other hand, monitoring the junction temperature is difficult and subject to inaccuracies (Valchev and Van den Bossche, 2005). Consequently, the main focus of this research is to track and forecast the evolution of $R_{DS_{ON}}$.

MOSFET end-of-life threshold criteria is not a deterministic value. Variations of $R_{DS_{ON}}$ are monitored, $\Delta R_{DS_{ON}}$, and it is generally assumed that a threshold value of $\Delta R_{DS_{ON}} 5\%$ would lead to failure of the MOSFET. Without assuming a fixed failure threshold, the main objective of this research is to forecast the evolution of $R_{DS_{ON}}$, which can be used to monitor the health of the MOSFET and predict the RUL, given a predefined failure threshold.

2.1 State-space modelling

For the following, the on-resistance of the discrete MOSFET devices is considered under a pristine condition R_{init} ; this value varies according to manufacturing tolerances. The temperature is assumed to be constant and equal to the value recommended by the manufacturer; this eliminates the temperature effects in the experiments.

The MOSFET ageing process results in an increment in the on-resistance with respect to the pristine condition, $\Delta R_{DS_{ON}}$, as a result of the degradation process due to the die-attach failure. An empirical degradation model for $R_{DS_{ON}}$ corresponds to exponential growth as a function of the time (Alonso et al., 2019):

$$\Delta R_{DS_{ON}}(t) = \alpha(e^{\beta t} - 1) \quad (1)$$

where α and β are model parameters that can be static (estimated using e.g. least-squares) or estimated on-line.

The degradation model in Eq. (1) can be converted into a dynamic model by considering the state-space representation. Assuming α and β as constant parameters and calculating the derivative of $\Delta R_{DS_{ON}}$ results in the following:

$$\dot{R}_{DS_{ON}}(k+1) = R_{DS_{ON}}(k)(1+\beta) + \alpha\beta \quad (2)$$

Equation (2) can be discretized and written in canonical form. Alternatively, if resistance under pristine conditions is considered, R_{init} , the state-space form results in the following (Dusmez et al., 2016):

$$R_{DS_{ON}}(t) = \alpha e^{\beta t} + R_{init} \quad (3)$$

$$R_{DS_{ON}}(k+1) = R_{DS_{ON}}(k)(1+\beta) - R_{init}\beta \quad (4)$$

3 MOSFET Ageing Forecasting Models

In order to forecast the future failure evolution of the MOSFET for different prediction horizons, different modelling approaches have been designed and tested: (i) classical state-space tracking algorithms based on Kalman filter variants, (ii) classical statistical forecasting methods, (iii) NN-based forecasting methods, and finally (iv) transformer-based forecasting methods.

3.1 Non-linear Kalman Filters

The Kalman filter can be used to predict the next step in linear processes, subject to process and measurement uncertainty. However, the state-space models (cf. Section 2.1) describe a non-linear process, and accordingly, non-linear variants of the KF are implemented. The [Extended Kalman Filter](#) (EKF) linearised the state-space equations around the operation point. The [Unscented Kalman Filter](#) (UKF) uses transformations to capture the propagation of the statistical properties of state estimates through the nonlinear equations. Thus, EKF and UKF can capture the effects of nonlinearities in the prediction of the aging process.

3.1.1 Extended Kalman Filtering (EKF)

The Extended Kalman Filter (EKF) (Chui and Chen, 2017) considers the estimation of the states $x \in \Re^n$ of discrete-time controlled processes that is non-linear. The EKF is a KF that linearizes the system equations around the mean using simple Taylor series with partial derivatives. The equations of the nonlinear stochastic difference equation for the state vector $x \in \Re^n$ are:

$$x_k = f(x_{k-1}, u_{k-1}, w_{k-1}) \quad (5)$$

with the measurement $z \in \Re^n$ is:

$$z_k = h(x_k, v_k) \quad (6)$$

the random variable w_k and v_k are the process and measurement noise, respectively ($w_k=0.002$ and $v_k=0.01$, in this work).

The estimation process for the non-linear systems begins with the linearisation of the system equations around the current operation point:

$$\begin{aligned} x_k &\approx \tilde{x}_k + A(x_{k-1} - \hat{x}_{k-1}) + Ww_{k-1} \\ z_k &\approx \tilde{z}_k + H(x_k - \tilde{x}_k) + Vv_k \end{aligned} \quad (7)$$

with x_k and z_k the actual state and the measurements and \tilde{x}_k and \tilde{z}_k the approximate state and measurement vector, respectively. The variable \hat{x}_k is the *a-posteriori* estimation of the state at step k . The approximate state and measurement vectors are calculated as follows:

$$\begin{aligned}\tilde{x}_k &= f(\hat{x}_{k-1}, u_{k-1}, 0) \\ \tilde{z}_k &= h(\tilde{x}_k, 0)\end{aligned}\tag{8}$$

The matrices \mathbf{A} and \mathbf{W} result from applying the Jacobian to the nonlinear functions:

$$\begin{aligned}A_{[i,j]} &= \frac{\partial f_{[i]}}{\partial x_{[j]}}(\hat{x}_{k-1}, u_k, 0) \\ W_{[i,j]} &= \frac{\partial f_{[i]}}{\partial w_{[j]}}(\hat{x}_{k-1}, u_k, 0) \\ H_{[i,j]} &= \frac{\partial h_{[i]}}{\partial x_{[j]}}(\tilde{x}_k, 0) \\ V_{[i,j]} &= \frac{\partial h_{[i]}}{\partial v_{[j]}}(\tilde{x}_k, 0)\end{aligned}\tag{9}$$

For simplicity in the notation, the subscript k in the Jacobians was omitted. The prediction error is defined as:

$$\tilde{e}_{x_k} = x_k - \tilde{x}_k\tag{10}$$

and the measurement residual:

$$\tilde{e}_{z_k} = z_k - \tilde{z}_k\tag{11}$$

The equations for the processes can be written as follows:

$$\begin{aligned}\tilde{e}_{x_k} &= A(x_{k-1} - \hat{x}_{k-1}) + \epsilon_k \\ \tilde{e}_{z_k} &= H(x_{k-1} - \tilde{x}_{k-1}) + \eta_k\end{aligned}\tag{12}$$

with ϵ and η are random variables with zero mean and covariance matrices WQW^T and VRV^T ; Q and R are covariance matrices of the noise and measurement noise with $p(w) \sim N(0, Q)$ and $p(v) \sim N(0, R)$, respectively. These equations resemble the differential and measurement equations of the discrete linear Kalman filter. Using superscript minus notation for the a priori estimation and in this occasion maintaining the discrete time subscript k for the Jacobians, the complete set of the EKF equations starts with the initialisation process with the initial estimate for $\hat{x}_k^- = E[x_k]$ and $P_k^- = E[(x_k - \hat{x}_k)(x_k - \hat{x}_k)^T]$. The time update process consists of the following processes:

- A. Project the state ahead:

$$\hat{x}_k = f(\hat{x}_{k-1}, u_{k-1}, 0)\tag{13}$$

- B. Project the error covariance ahead:

$$P_k^- = A_k P_{k-1} A_k^T + W_k Q_{k-1} W_k^T\tag{14}$$

The measurement update consists of the following steps:

- A. Compute the Kalman gain:

$$K_k = P_k^- H_k^T (H_k P_k^- H_k^T + V_k R_k V_k^T)^{-1}\tag{15}$$

- B. Update estimates with measurements z_k :

$$\hat{x}_k = \hat{x}_k^- + K_k (z_k - h(\hat{x}_k^-, 0))\tag{16}$$

- C. Update the covariance error:

$$P_k = (I - K_k H_k) P_k^-\tag{17}$$

The procedure is applied to the preceding equations that model the MOSFET degradation process Eqs. (2)-(4) in order to estimate the values of the describing parameters α and β .

3.1.2 Unscented Kalman Filtering (UKF)

The UKF is centred around the idea to describe the state x as a distribution defined by a small number of characteristic sampling points, named sigma points, that represent the Gaussian underlying distribution. The UKF operates in an iterative way through sigma point determination, state-estimation, and update stages.

Sigma point determination

The state-estimation, x_k is approximated by the by the sigma points, which are calculated from the mean, \bar{x}_k , and covariance, P_k of the previous step. For each iteration, k , the sigma points for the state variables are estimated as follows (Wan and van der Merwe, 2001):

$$\begin{aligned} x[0] &= \bar{x}_k \\ x[i] &= \bar{x}_k + (\sqrt{(n + \lambda)P_k})_i; \text{ for } i=1,\dots,n \\ x[i] &= \bar{x}_k - (\sqrt{(n + \lambda)P_k})_{i-n}; \text{ for } i=n+1,\dots,2n \end{aligned} \quad (18)$$

where λ is a scaling parameter, $\lambda = \alpha^2(n + k) - n$, and α , β and K are all tunable parameters to determine the spread of the sigma points.

These sigma points are then fed into the state equation to generate a set of new sigma points for the state variables in the current step. After calculation of sigma points, they are passed through the state-estimation function in Eq. (1), and their weights are calculated. Finally, they are used to form a state prediction together with a covariance and cross-correlation matrix.

Update stage

In the update step, the predictions are combined with measurements [cf. Eq. (2)] to form the new estimate. A detailed description of the UKF can be found in (Wan and van der Merwe, 2001).

3.2 Classical Statistical Forecasting

3.2.1 Autoregressive Integrated Moving Average (ARIMA)

ARIMA combines differencing with autoregression and moving average. The model can be written as follows (Hyndman and Athanasopoulos, 2018):

$$\hat{y}_t = c + \phi_1 y'_{t-1} + \dots + \phi_p y'_{t-p} + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q} + \epsilon_t \quad (19)$$

where \hat{y}_t is the predicted value of the time-series. The precursors of the time series include lagged values of y_t and lagged errors. The constants to estimate include the p order of the autoregressive part, the d degree of the first differencing involved, and the q order of the moving average part.

These parameters have been adjusted with the Hyndman-Khandakar algorithm. This algorithm combines unit root tests, minimization of the Akaike Information Criteria (AIC) and maximum likelihood estimation to obtain an ARIMA model. The implementation of the ARIMA model is developed using the `forecasting` package of R (Hyndman et al., 2024).

3.2.2 Holt's Linear Trend Method

Applies exponential smoothing to univariate data, with trend and no seasonal pattern. [Holt's linear trend method](#) can be described with the following equations (Hyndman and Athanasopoulos, 2018):

$$\hat{y}_{t+h|t} = l_t + h b_t \quad (20)$$

$$l_t = \alpha^* y_t + (1 - \alpha^*)(l_{t-1} + b_{t-1}) \quad (21)$$

$$b_t = \beta^*(l_t - l_{t-1}) + (1 - \beta^*)b_{t-1} \quad (22)$$

where l_t is the smoothed level of the series at instant t , b_t denotes an estimate of the trend of the series at time t , α^* is a smoothing parameter for the level and β^* is the smoothing parameter for the trend.

The application of Holt's method requires smoothing parameters α^* , β^* , and initial values l_0 and b . All forecasts can be computed from the data after computing the initial values. Initial values are selected using the minimization of the sum of squared errors (SSE), defined as follows:

$$SSE = \sum_{t=1}^T (y_t - \hat{y}_{t|t-1})^2 \quad (23)$$

This formulation requires a nonlinear optimization. The implementation of the forecasting package in R (Hyndman et al., 2024) is adopted in this work.

3.3 Forecasting with Neural Networks (NN)

NN models can operate as universal approximators of nonlinear functions by learning complex relationships between the response variable and its predictors (Scardapane, 2024). NN models are made up of an input layer, hidden layers, and an output layer. The input for time-series forecasting is comprised of historical time-series signals up to a input window size p , $\mathbf{x} = \{x_t, x_{t-1}, \dots, x_{t-p+1}\}$. The hidden layer is defined as follows:

$$h^{(l)} = \sigma(\mathbf{W}^{(l)} h^{(l-1)} + b^{(l)}) \quad (24)$$

where $h^{(l)}$ is the output of the l -th layer (for the first layer $h^{(0)} = \mathbf{x}$), $\mathbf{W}^{(l)}$ is the weight matrix of the l -th layer, $b^{(l)}$ is a bias vector of the l -th layer, and $\sigma(\cdot)$ is a nonlinear activation function.

Finally, the output layer produces a h steps ahead forecast, \hat{x}_{t+h} , chaining hidden layers and the output layer:

$$\hat{x}_{t+h} = \mathbf{W}^{(L)} h^{(L-1)} + b^{(L)} \quad (25)$$

where L is the total number of layers.

Forecasting with NNs requires the specification of a NN architecture that can approximate and extrapolate the underlying data-generating process. Namely, for time-series forecasting, it is necessary to evaluate the forecasting information of lagged signals and, if present, include seasonality. Accordingly, feature selection for time-series forecasting with NN needs to consider autoregressive and seasonality components.

In this context, to improve the accuracy of the prediction through a set of NN prediction models, (Kourentzes et al., 2014) developed NN ensemble operators for the prediction of time series, with excellent prediction results in different prediction competitions (Wang et al., 2023). This approach combines multiple base models in an ensemble model. Through the nnfor package in R, this work implements ensemble forecasting models based on [Neural Network](#) (E-NN) and [Extreme Learning Machines](#) (E-ELM) models (Kourentzes, 2023). ELMs mitigate the optimisation of NN. Namely, instead of attempting to tune all weights in a NN, they are left to their random initial values, except for weights in the output layer.

In all the E-NN and E-ELM model configurations, a single hidden layer is used (Kourentzes et al., 2014). Regarding the number of neurons, it was evaluated by five-fold CV, using configurations from 2 up to 20 neurons. First-level differences are calculated over the input signal and lagged signals are evaluated in predefined windows and best lagged signals are selected. The series is also evaluated for seasonality and, if it is stochastic, seasonal differences are added to the input of the model (Kourentzes et al., 2014). We use 20 ensemble members that are combined with the median operator. The activation function for E-NN and E-ELM is a sigmoid function. The output layer of the E-ELM is trained using Lasso regression. Backpropagation is used to train adjust the weights of each neuron by calculating the gradient of the loss function.

3.4 Temporal Fusion Transformers (TFTs)

The architecture of TFTs was introduced in (Lim et al., 2021). TFTs are particularly effective for time series forecasting due to (i) their increased ability to capture long- and short-term dependencies through the self-attention mechanism; (ii) adaptive attention mechanism for dynamic temporal processing to prioritize important time points and ignore irrelevant information; (iii) ability to handle multiple data types including past observations, known future inputs and static covariates; and (iv) interpretable results which offer insights into which variables and time periods are most important for making predictions.

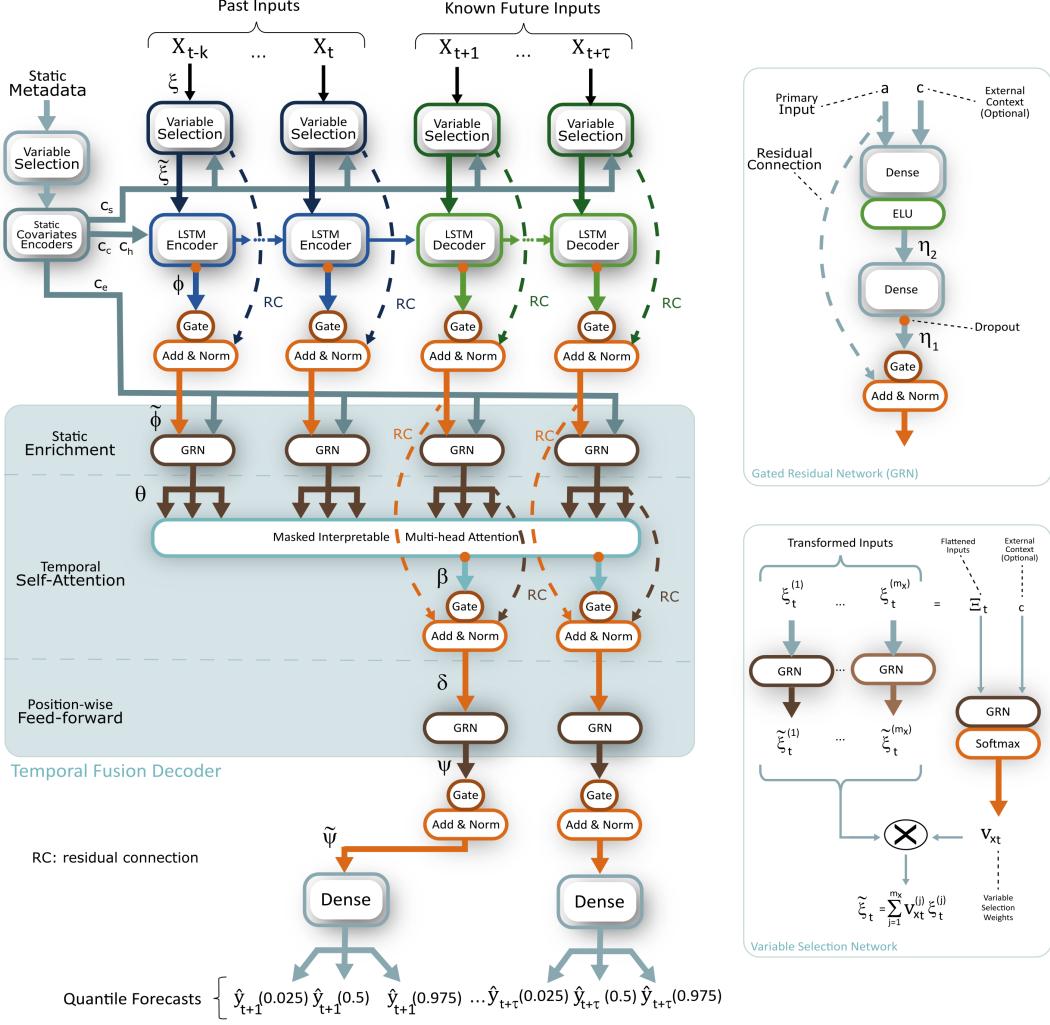


Figure 2: TFT architecture including Gated Residual Networks, Variable Selection Networks and encoder-decoder structure.

The TFT architecture is designed using canonical components to build feature representations for improved forecasting performance, which is shown in Figure 2. This includes:

- gating mechanism*, which skips unused components, provides adaptive depth and network complexity;
- variable selection networks* to select relevant variables at each time step;
- static covariate encoders* to integrate static features into the network through encoding context vectors to condition temporal dynamics;
- temporal processing* to learn long and short term temporal relationships from the inputs. A sequence-to-sequence LSTM layer is employed for local processing and long-term dependencies are captured through a multi-head attention block;
- Finally, *prediction intervals* via quantile forecasts to determine the range of likely values for each prediction horizon.

Gating mechanisms are used to regulate the contribution of a sample, a , to a context vector c through Gated Residual Networks (GRN), defined as follows:

$$GRN_w(a, c) = \text{Layer Norm}(a + \text{GLU}_w(\eta_1)) \quad (26)$$

$$\eta_1 = W_{1,w}\eta_2 + b_{1,w} \quad (27)$$

$$\eta_2 = ELU(W_{2,w}a + W_{3,w}c + b_{2,w}) \quad (28)$$

where ELU is the Exponential Linear Unit activation function, $\eta_1 \in R^{d_{model}}$ and $\eta_2 \in R^{d_{model}}$ are intermediate layers, LayerNorm is standard layer normalization, and w is an index to denote weight sharing. The parameter a represent the linear contribution. Component gating layers based on gated linear units (GLUs) are adopted to control the extent of nonlinear contribution to suppress unnecessary information in a given data set. The GLU is described as follows:

$$GLU_w(\gamma) = \sigma(W_{4,w}\gamma + b_{4,w}) \odot (W_{5,w}\gamma + b_{5,w}) \quad (29)$$

where γ is the input, σ is the sigmoid activation function, $W_{(.)} \in R^{d_{model} \times d_{model}}$ and $b_{(.)} \in R^{d_{model}}$ are the weights and biases, \odot is the element-wise Hadamard product, d_{model} is the hidden state size.

The variable selection networks (VSN) ignore noisy inputs and concentrate on valuable data to improve model performance. VSN select the most relevant features from the input time-series data according to the forecasting horizon through a weighing mechanism. The variable selection weights provide interpretability of input variable for the TFT model. Let $\xi_t^{(j)}$ denote the transformed input of the j -th variable at time t , with being the flattened vector $\Xi_t = [\xi_t^{(1)^T}, \dots, \xi_t^{(m_x)^T}]^T$ of all past inputs at time t . The flattened vectors and the static covariate c_s are fed into a GRN and pass through a Softmax layer to calculate the selection weight of the variable.

$$v_{X_t} = \text{Softmax}(GRN_{V_x}(\Xi_t, c_s)) \quad (30)$$

$$\xi_t^j = GRN_{\xi_t^j}(\xi_t^{(j)}) \quad (31)$$

$$xi_t = \sum_{j=1}^{m_x} v_{x_t}^{(j)} \xi_t^{(j)} \quad (32)$$

where $v_{X_t}^{(j)}$ is the j -th element of the vector v_{x_t} . The variable selection weight provides the explanatory properties for the forecasting of the TFT model, *i.e.* the globally significant variables of the forecasting task can be identified in Eq. (31), each is non-linearly processed by its GRN. Finally, the processed features are weighted by their variable selection weights and combined as shown in Eq. (32).

The attention function can be summarized as mapping three vectors, a query and a key-value pair, to an output that is expressed as follows:

$$\text{Attention}(Q, V, K) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (33)$$

where d_k is the dimension of Q , $Q \in R^{n \times d_k}$ is the query, $K \in R^{m \times d_k}$, and $V \in R^{m \times d_v}$ are the key-value data pair. The attention function can be seen as a mapping sequence from $Q \in R^{n \times d_k}$ to $Q \in R^{n \times d_v}$ by a dot product of three matrices with sizes $n \times d_k$, $d_k \times m$ and $m \times d_v$.

Self-attention is a commonly used attention mechanism to indicate the attention distribution of a time-sequence. Self-attention can be defined as $\text{Attention}(X, X, X)$, where $X \in R^{n \times d_k}$. Performing several attention functions in parallel with different projections is more beneficial than using a single attention function. Accordingly, multi-head attention is adopted to jointly yield the outputs from different mapping subspaces, which is expressed as:

$$\text{Multihead}(Q, V, K) = [H_1, \dots, H_{m_H}]W_H \quad (34)$$

$$H_h = \text{Attention}(QW_Q^{(h)}, KW_K^{(h)}, VW_V^{(h)}) \quad (35)$$

where $W_Q^{(h)}$, and $W_K^{(h)} \in R^{d_{model} \times d_{attn}}$, $W_V^{(h)} \in R^{d_{model} \times d_v}$, are head-specific weights for the keys, queries and values, and $W_H \in R^{d_{model} \times (d_v m_H)}$ linearly combines outputs concatenated from all the heads, H_h .

Given that each head uses a different value, the individual attention weight cannot represent the importance of a particular feature. The multi-head attention is thus changed to the shared value in each head:

$$\text{Interpretable Multi Head}(Q, K, V) = \tilde{H} W_H \quad (36)$$

$$\tilde{H} = \tilde{A}(Q, K) V W_V \quad (37)$$

$$\tilde{H} = \frac{1}{m_H} \sum_{h=1}^{m_H} \text{Attention}(Q W_Q^{(h)}, K W_K^{(h)}, V W_V^{(h)}) \quad (38)$$

The overall TFT architecture is framed with an LSTM encoder-decoder to capture the short-term memory, and the multi-head attention captures the long-term dependencies. The temporal fusion decoder, uses a series of (i) locality enhancement with sequence-to-sequence layers, (ii) static enrichment layer, (iii) temporal self-attention layer, and (iv) position-wise feedforward layer; to learn temporal relationships present in the dataset (Lim et al., 2021).

The encoder-decoder structure enables the static metadata to influence local processing. This is formalised through context vectors, c_c, c_h , static covariate encoder context c_e , and an external context vector c_s . Gated skip connection layers are applied and a static enrichment layer that enhances temporal features with static metadata.

After this stage, the self-attention is applied through static enriched temporal features and interpretable multi-head attention. Additionally, a gating layer is applied to facilitate training, which results in $\delta(t, n)$. Subsequently, a non-linear processing is applied to the outputs of the self-attention layer:

$$\psi(t, n) = GRN_\psi(\delta(t, n)) \quad (39)$$

As shown in Figure 2, a gated residual connection is also applied, which results in $\tilde{\psi}$. Finally, quantile forecasts are generated using a linear transformation of the output of the temporal fusion decoder, $\tilde{\psi}(t, \tau)$, as follows:

$$\hat{y}(q, t, \tau) = W_q \tilde{\psi}(t, \tau) + b_q \quad (40)$$

where $W_q \in \mathbb{R}^{1 \times d}$, $b_q \in \mathbb{R}$ are linear coefficients of the quantile q .

Finally, the TFT is trained by jointly minimizing the quantile loss function, summed across all quantile outputs:

$$\mathcal{L} = \sum_{y_t \in \Omega} \sum_{q \in Q} \sum_{\tau=1}^{\tau_{max}} \frac{QL(y_t, \hat{y}(q, t - \tau, \tau), q)}{M \tau_{max}} \quad (41)$$

$$QL(y, \hat{y}, q) = q(y - \hat{y})_m + (1 - q)(\hat{y} - y)_m \quad (42)$$

where Ω is the training domain, which includes M samples, W represents TFT weights, Q is the set of quantiles ($Q = \{0.025, 0.5, 0.975\}$ in this work), and $(.)_m = \max(0, .)$.

Transformers are known to perform very well on long-term time series (Biggio et al., 2023). In this research, the TFT package Darts available for Python is used (Herzen et al., 2022). Overall, TFT shows very good generalization performance to unseen degradation conditions, and it is able to handle arbitrarily complex ageing profiles. This is the case even when these are chosen significantly outside the training distribution. There are several hyperparameters that need to be considered when using the TFT. Consequently, different architectures were tested to select the best combination of three of those hyperparameters, including hidden size, LSTM layers and attention heads, as shown in Table 2.

In order to determine the best TFT architecture for different prediction horizons, the performance of LSTM layers, number of neurons, and attention heads were tuned.

Table 2: Tested TFT hyperparameters.

Model name	Hidden size	LSTM layers	Attention heads
A	128	4	4
B	128	8	4
C	128	4	2
D	256	4	4
E	256	8	4
F	256	4	2
G	64	1	4

4 Case Study

The MOSFET dataset provided by NASA (Celaya et al., 2012) is used to evaluate forecasting ability of different models in different configurations, time-scales, and learning proportions.

4.1 Experimental Environment

Accelerated aging experiments are used to assess reliability in a much shorter time than long-term reliability tests. Specifically, thermal cycles are applied to create thermomechanical stresses in the switching devices due to the mismatch in the thermal expansion coefficients of different package elements. Failure conditions include latch-up, thermal runaway, or failure to turn on due to gate damage. In thermal cycles, no heatsink is used, significantly reducing heat dissipation, allowing for self-heating during commutes. The case temperature of the semiconductor device was measured using a thermocouple attached to the copper case flange (drain).

In this case, the accelerated thermal stress procedures were implemented for IRF520NPbF MOSFETs in run-to-failure experiments. For thermal cycling, on/off cycles were produced by applying a gate voltage of 15 V with a PWM signal at a frequency of 1 kHz and a duty cycle of 40%. The drain-source was biased at 4 V DC, and a resistive load of 0.2 Ω was used on the source output. The dataset included gate-source voltage (V_{GS}), drain-source voltage (V_{DS}), drain current (I_D), and flange (case) temperature (T_f). The sampling rate of these time series is 1 MHz, with a transient response measurement every 400 ns.

X-ray and scanning acoustic images after degradation show that the die-attach solder migrates, forming voids. Die-attach damage due to thermal cycling reduces the contact area between solder-copper and solder-silicon interfaces, decreasing heat flow compared to a pristine die-attach. This results in a higher junction temperature and an increased ON resistance for the degraded device. Hence, the drain-to-source resistance ($R_{DS_{ON}}$) can be considered a precursor to assess the state of health (Celaya et al., 2012). Furthermore, junction temperature is also a function of the case temperature, and the measured $R_{DS_{ON}}$ was normalized to eliminate the case temperature effects and reflect only changes due to degradation.

Due to manufacturing variability, not all transistors show an identical starting value of $R_{DS_{ON_0}}$ at time $t = 0$, which is defined as pristine condition. For fair comparison, the initial $R_{DS_{ON_0}}$ value (“pristine condition value”) is subtracted from the absolute measurements. The normalized time series results in a trajectory ($\Delta R_{DS_{ON}}$) from pristine condition to failure, representing the degradation process due to die-attach failure and the increase in $R_{DS_{ON}}$ through the aging process.

Therefore, for prognostics, it is assumed that die-attach failure is the only active degradation mechanism during the accelerated aging experiment. The single value of $\Delta R_{DS_{ON}}$ is used to assess the health state of the device, with a failure threshold of a 0.05 increase in $\Delta R_{DS_{ON}}$ (Celaya et al., 2012).

4.2 Pre-processing Experiments

The main pre-procesing stages are divided into data sampling, temperature filtering, resistance normalization and temporal filtering. Namely, during the data sampling stage, firstly, the gate voltage V_{GS} is used to separate ON and OFF states of the MOSFET. Figure 1 shows the ON and OFF cycles of one experiment. The on-resistance $R_{DS_{ON}}$, which correlates with the junction temperature, is computed as the ratio of V_{DS} to I_D during the on-state (with duration T_{ON}) of the square waveform (Celaya et al., 2012):

$$R_{DS_{ON}}(k) = \sum_{i=1}^{T_{ON_i}} \frac{V_{DS}(i)}{I_D(i)} \quad (43)$$

where $k = \{T_{ON_1}, \dots, T_{ON_K}\}$ are the different ON cycles during the MOSFET accelerated thermal ageing process, and $i = \{1, \dots, T_{ON_i}\}$ is the duration of i-th ON cycle of the MOSFET.

Subsequently, $R_{DS_{ON}}$ values computed at flange temperature values below the $T_f < \text{Low Temp}$ threshold were removed from the data in order to discard $R_{DS_{ON}}$ values with no influence on thermal ageing. Furthermore, in order to eliminate the case temperature effects and observe only the changes produced by degradation, resistance normalization was carried out subtracting the pristine condition to all the $R_{DS_{ON}}$ values (Celaya et al., 2012). Finally, minute-based mean values of $\Delta R_{DS_{ON}}$ values were computed to filter out noise.

Figure 3 shows the calculation example of $R_{DS_{ON}}$, including V_{DS} , V_{GS} and I_D , that are used to compute $R_{DS_{ON}}$ [cf. Eq. (43)].

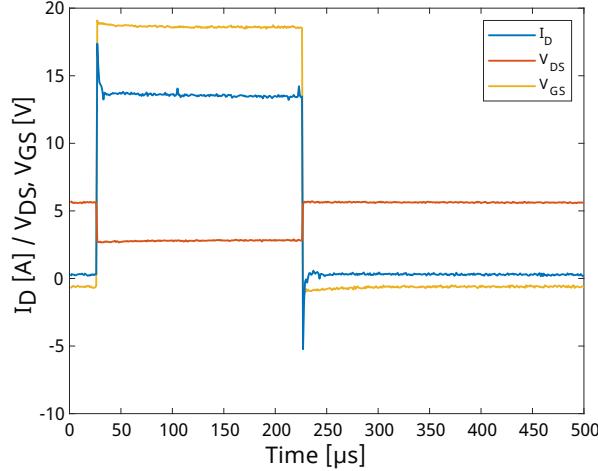


Figure 3: V_{DS} , V_{GS} and I_D used to compute $R_{DS_{ON}}$.

After this process, four experiments were selected from the total of 42 experiments. This selection process was done based on the amount of remaining data points for each experiment after completing the pre-processing stage and evaluating the instant when failure condition was reached. Accordingly, Figure 4 shows the available run-to-failure experiments of the MOSFETs after the pre-processing stage.

Figure 4 shows that there is an inherent variability in the ageing processes of the same MOSFETs for different tests. This variability may have implications when designing an ageing forecasting model for MOSFETs of the same family.

4.3 Statistical properties & Evaluation Metrics

All the available experiments are non-stationary, which have been evaluated through the Dickey-Fuller test. In addition, through the analysis of the autocorrelation, it has been observed that all the experiments have a trend and they do not have seasonality. These characteristics make the experimental datasets good candidates to apply ARIMA and Holt models.

The evaluation metric Mean Average Percentage Error (MAPE) is defined as follows:

$$MAPE = \frac{100}{N} \sum_{t=1}^{t=N} \left| \frac{y_t - \hat{y}_t}{y_t} \right| \quad (44)$$

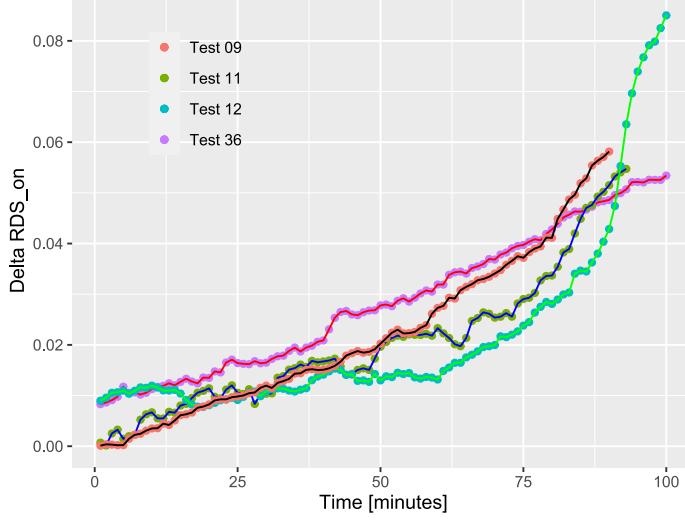


Figure 4: Obtained run-to-failure experiments after pre-processing.

5 Numerical Experiments

In this section, MOSFET ageing forecasting models have been designed and tested for short- and long-term predictions.

5.1 Short-term MOSFET Ageing Forecasts

5.1.1 Model design

The performance of the selected model architectures was evaluated by training and testing in an iterative way. Given a prediction horizon of N steps ahead, the process starts training with the first p samples of a given experiment and makes a 1 step-ahead prediction. If $N > 1$, then, forecast results are concatenated with the training set and thus, the training set is incremented in N . Then again, another N step ahead prediction is performed. This process is repeated until the prediction of the last available point of the experiment is carried out. Algorithm 1 describes the details of the training-testing process. Initial tests have been performed using 33 % of the available dataset for each experiment.

For each of the forecasting models (EKF, UKF, ARIMA, Holt, TFT, E-ELM, E-NN), in Line 2 of the Algorithm 1 it is necessary to adjust the hyperparameters to fit the model. In the TFT model, a hyperparameter tuning process is needed to select to adjust the parameters (see Table 2). The obtained results are presented in Figure 5.

[Figure 5 shows that there is a direct dependence between the hyperparameters and obtained results.](#) The best TFT architecture was selected by counting the number of times that each architecture performed best among all solutions. [In this case, the model G performed best among all configurations \(cf. Table 2\).](#) The hyperparameter tuning process is repeated for all models with their respective parameters. After fitting the model, predictions are made iteratively across the validation set. Once the predictions are obtained, \mathcal{V} , the MAPE of the prediction results is computed.

5.1.2 Prediction Results

Table 3 displays the obtained mean MAPE results for MOSFET ageing predictions at different forecast horizons $N = \{1, 2, 4\}$ for different experiments and different predictive models using the available 33% of the datasets as training sets.

Table 3 shows that, consistently across different forecasting horizons, E-ELM performs better for tests #9 and #36, except N=4 for test #9 with HOLT's model, which in practice is also a filter based on weighted past values, and for the tests #11 and #12, ARIMA performs best.

The E-ELM handles well tests #9 and #36 because it lacks of underlying model-form and noise hypotheses. It models the future evolution from a purely black-box perspective, including neuron weights and biases optimized to mimic the $R_{DS_{ON}}$ evolution, based on the selected delayed signals. In contrast, it can be observed that the tests #11 and #12

Algorithm 1 Short-term training.

Inputs:

- $D = [d_1, \dots, d_K]$ ▷ Data vector of K elements.
- p ▷ Data-points to fit the model, $p < K$.
- $g(.)$ ▷ Forecasting model.
- $D_T = [d_1, \dots, d_p]$ ▷ Training vector, p elements.

Output: \mathcal{V}

- $h = K - p$ ▷ Predictions vector.
- 1: **for** $(i = 1 : h)$ **do** ▷ Prediction horizon.
- 2: $g = fit(D_T)$ ▷ For data length.
- 3: $\mathcal{F} = g(D_T, 1)$ ▷ Fit model g with p data-points.
- 4: $\mathcal{V} \leftarrow \mathcal{F}$ ▷ 1-step ahead forecast.
- 5: **if** $N > 1$ **then** ▷ Save forecast in V
- 6: **for** $(j = 2 : N)$ **do** ▷ Multi-step ahead forecast.
- 7: $D_T = D_T \cup \mathcal{F}$ ▷ Concatenate.
- 8: $g = fit(D_T)$ ▷ Re-fit model g .
- 9: $\mathcal{F} = g(D_T, 1)$ ▷ 1-step ahead forecast.
- 10: $\mathcal{V} \leftarrow \mathcal{F}$ ▷ Save forecast in V
- 11: **end for**
- 12: **end if** ▷ Increase p in N
- 13: $p = p + N$ ▷ Update training vector.
- 14: $D_T = [d_1 : d_p]$
- 15: **end for**
- 16: **return**(\mathcal{V})

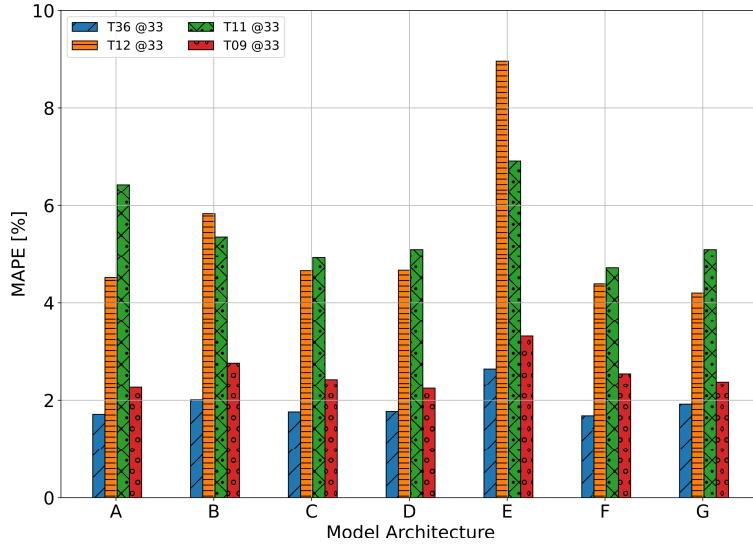


Figure 5: TFT architecture selection for short-term MOSFET ageing forecasting for different tests (T09, T11, T12, T36) with 33% of data used for training.

include more variability. Accordingly, this reflects the capability of the ARIMA model to deal with variable signals, filtering out local fast variations that have no relationship with long-term behaviour.

Among state-space models, UKF is the most accurate for all tests and different prediction horizons. Although it does not obtain the most accurate results for the tests #9 and #36, UKF results are comparable to the rest of the models. In contrast, for the EKF model with tests #11 and #12 the accuracy drops substantially. Clearly in this case, the Taylor series approximation of the EKF was not sufficient to capture the effects the non-linearities present in the model. The EKF lacks the ability to adapt to changes in signal modifications, resulting in continuous oscillations of the state-space parameters α and β , which have a direct impact on the error. In contrast, UKF results are better because the parameters are more stable and adapt to signal modifications. Figure 6 shows the influence of the modification of α and β [parameters in the EKF and UKF models](#) with respect to the learning horizon.

Table 3: Mean $\Delta R_{DS_{ON}}$ prediction MAPE for different tests, models and steps-ahead (N) — best results highlighted.

Test	Model	Forecasting Horizon		
		N=1	N=2	N=4
#36	EKF	2.61	3.65	5.44
	UKF	2.09	2.82	4.58
	HOLT	1.65	1.84	1.93
	ARIMA	1.6	1.93	2.17
	TFT	1.92	2.51	3.39
	E-NN	1.9	2.2	3.1
#9	E-ELM	1.6	1.8	1.9
	EKF	3.39	4.47	7.37
	UKF	3.31	4.67	7.45
	HOLT	2.19	2.55	3.47
	ARIMA	2.04	2.5	3.56
	TFT	2.37	3.35	5.33
#11	E-NN	2.3	3	5
	E-ELM	2	2.5	3.9
	EKF	7.12	11.1	14.3
	UKF	4.38	6.39	8.44
	HOLT	4.51	5.60	7.80
	ARIMA	3.86	5.48	7.57
#12	TFT	5.09	6.73	10.5
	E-NN	5.1	6.8	9.8
	E-ELM	4.5	6.1	8.8
	EKF	8.90	13.7	18.2
	UKF	5.54	6.38	8.89
	HOLT	4.57	5.11	8.17
#12	ARIMA	3.68	5.02	7.15
	TFT	4.20	7.03	11.1
	E-NN	5	6.8	9.8
	E-ELM	4.5	6.1	8.8

Figure 6 shows that the UKF model parameters adapt and stabilize as the prediction time progresses, whereas the EKF model parameters are non-stable, aligned with the limitations of the EKF to match the nonlinear characteristics outside the operation point, and this is reflected in the increased prediction error.

As for the TFT model for short-term predictions, it can be observed that the performance is in between the best and worst performing models. Increasing the forecasting horizon increases the MAPE of all models across different tests.

5.1.3 Influence of train-test proportions

In order to evaluate the influence of different training-testing proportions on the forecasting accuracy, along with the TFT model, the best performing state-space statistical forecasting and NN with AR inputs models have been tested. They correspond to *i.e.* UKF, ARIMA and E-ELM models for different tests and different forecasting horizons. Figure 7 shows the obtained results.

Figure 7 shows that, consistently across all models, the more training data, the lower the MAPE. Moreover, the longer the prediction horizon, the greater the MAPE. There are some exceptions, especially for the UKF model, where the training set determines the initial conditions. Namely, the corresponding α and β values are calculated through least squares using the training dataset, and then these parameters are used as initial values for the state-estimation algorithms. This is not always an information gain for the state-space model based forecasting, as observed in test #12.

As for the ARIMA and E-ELM, parameters of the trained models fully depend on the available data. Therefore, if longer datasets include the behaviour of the testing set, then the accuracy and uncertainty should be lowered. However, if the testing set includes out-of-distribution samples, the ability to model the ageing will become smaller.

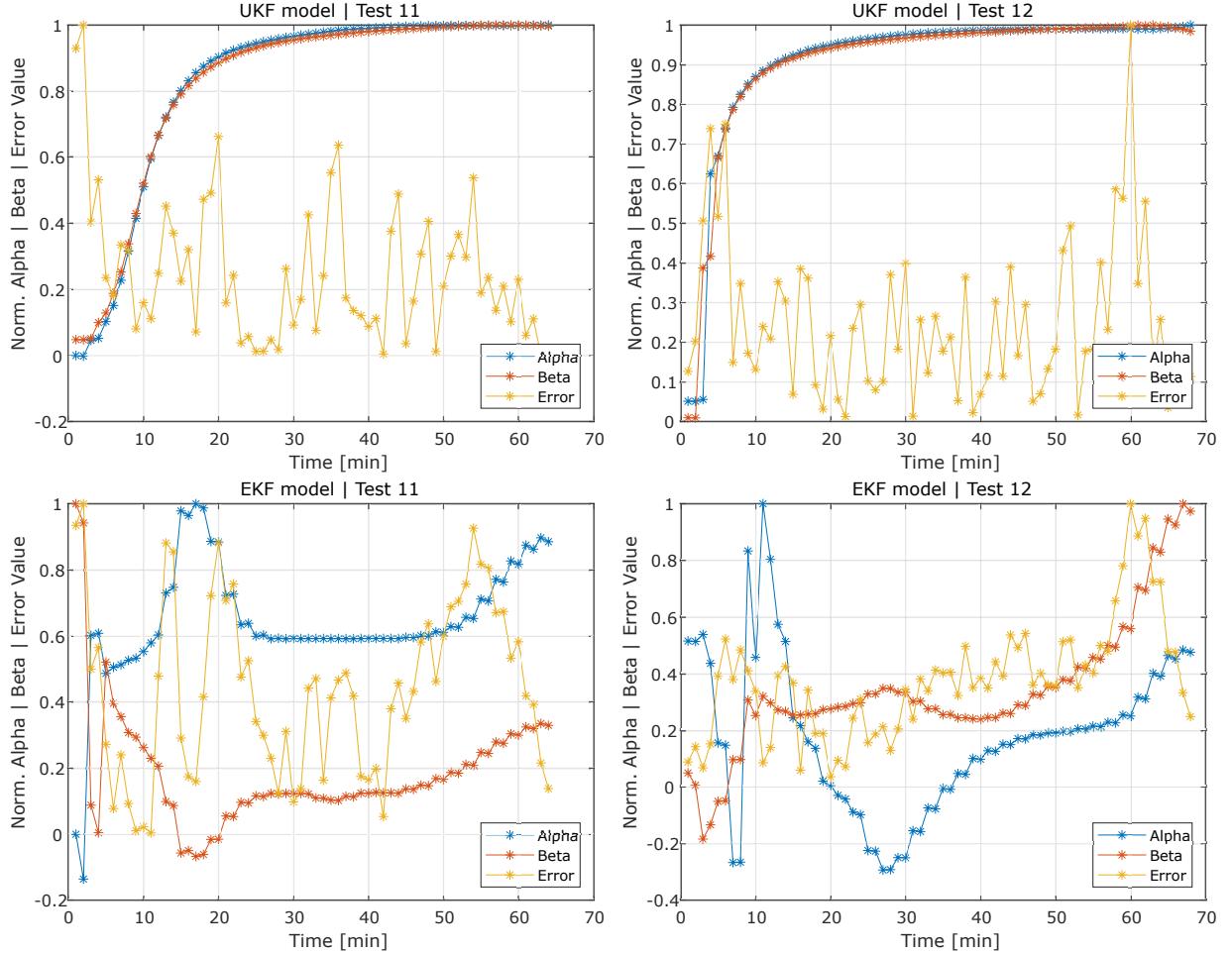


Figure 6: Adaptation of α and β parameters with respect to prediction error for $N=2$ for the tests #11 and #12 for the UKF and EKF models.

5.1.4 Influence of the state-transition model

It is possible to test the state-space models with alternative definitions. Namely, Table 4 shows the mean MAPE results for the state transition model in Eq. (4) using 33% of the datasets for training.

Table 4: Mean $\Delta R_{DS_{ON}}$ prediction MAPE for the UKF model.

Test	Steps-ahead		
	N=1	N=2	N=4
#9	3.49	5.02	8.12
#11	4.76	7.18	10.00
#12	4.97	6.70	9.99
#36	2.80	4.15	7.19

Comparing the UKF model results displayed in Table 3 with Table 4, which uses the state-space transition model in Eq. (2), it can be observed that the MAPE in the latter case is greater for all the different tests (except, test #12, N=1), demonstrating the superior capabilities of the state-space model in Eq. (2) for modelling non-linear processes.

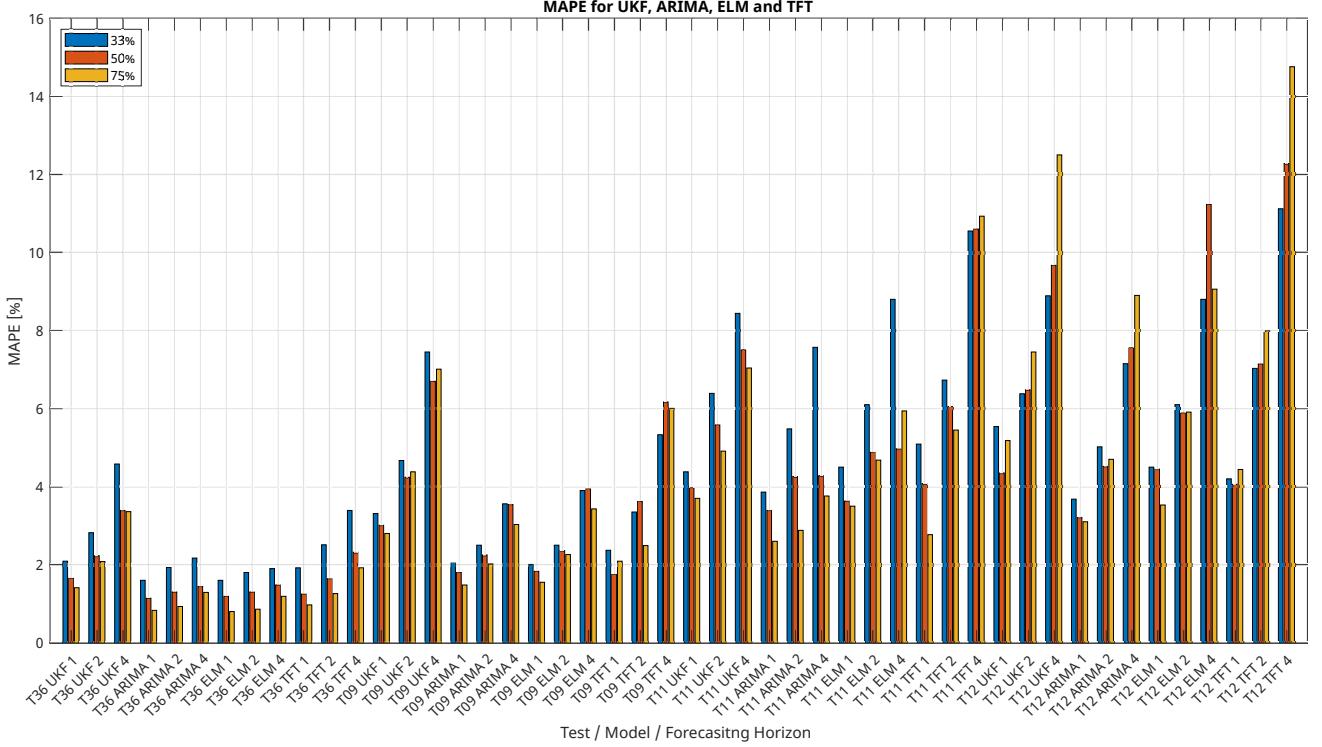


Figure 7: MAPE values for different training proportions, tests and models.

5.2 Long-term MOSFET Ageing Forecasting

In order to evaluate the generalisation capability of the designed models on unseen experiments, long-term forecasting models have been designed and tested. To this end, the TFT model architecture has been used, which has shown good results in various long-term prediction applications (Biggio et al., 2023). In addition, the best model from Table 3 has been selected for comparison purposes, *i.e.* the ELM model.

Based on the assumption that a Leave One Out (LOO) testing is designed for MOSFETs of the same family, firstly, outlier experiments are discarded, which indicates different MOSFET technology. To this end, the Fréchet distance was calculated among different experiments (Har-Peled and Raichel, 2014). Table 5 displays the obtained results.

Table 5: Similarity of experiments based on Fréchet distance.

Test	#9	#11	#12	#36
#9	0	0.003	0.027	0.008
#11	0.003	0	0.03	0.008
#12	0.027	0.03	0	0.032
#36	0.008	0.008	0.032	0

It is observed that the distance of test #12 with respect to other tests is greater. This potentially indicates that it pertains to a different MOSFET family, and therefore, it is not considered for LOO experiments.

5.2.1 Model design

One of the strengths of the TFT is that it can incorporate covariates that impact on model predictions. Accordingly, to inform the long-term performance of the MOSFET aging, a covariate has been designed based on the piecewise linearisation of the training tests.

The piecewise linearisation consists of (i) splitting the tests into blocks of 20% data samples (ii) curve-fitting through least squares for each block and (iii) finally, the covariate is defined as a piecewise function of the fitted curves. This is described in Eq. (45), where $\{m_1, \dots, m_5\}$ and $\{b_1, \dots, b_5\}$, correspond to the linear curve parameters for the fitted curve.

$$y(t) = \begin{cases} m_1 t + b_1 & 0 \leq t \leq 19 \\ m_2 t + b_2 & 20 \leq t \leq 39 \\ m_3 t + b_3 & 40 \leq t \leq 59 \\ m_4 t + b_4 & 60 \leq t \leq 79 \\ m_5 t + b_5 & 80 \leq t \leq 99 \end{cases} \quad (45)$$

Figure 8 shows an example of the covariates used for test 36 through the piecewise linearization of tests 9 and 11.

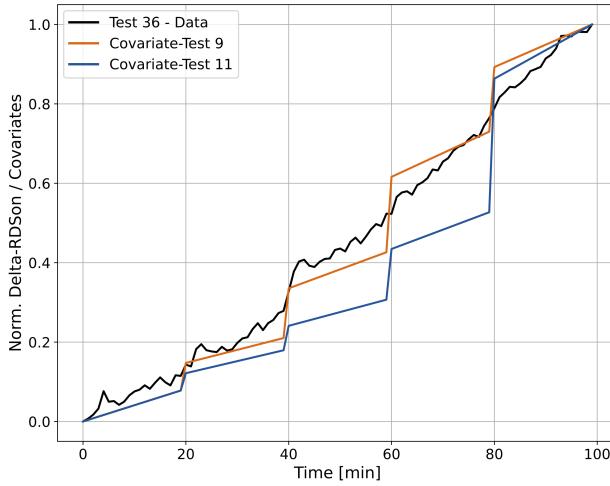


Figure 8: Normalized test 36 with the designed covariates using tests 9 and 11.

To select the long-term architecture, different combinations of hyperparameters have been tested (cf. Table 2). Figure 9 shows the results of the selection of the model architecture for different tests (T09, T11, T12, T36), with different training data proportions (30, 50, 70), and including TFT models without and with covariates (wC).

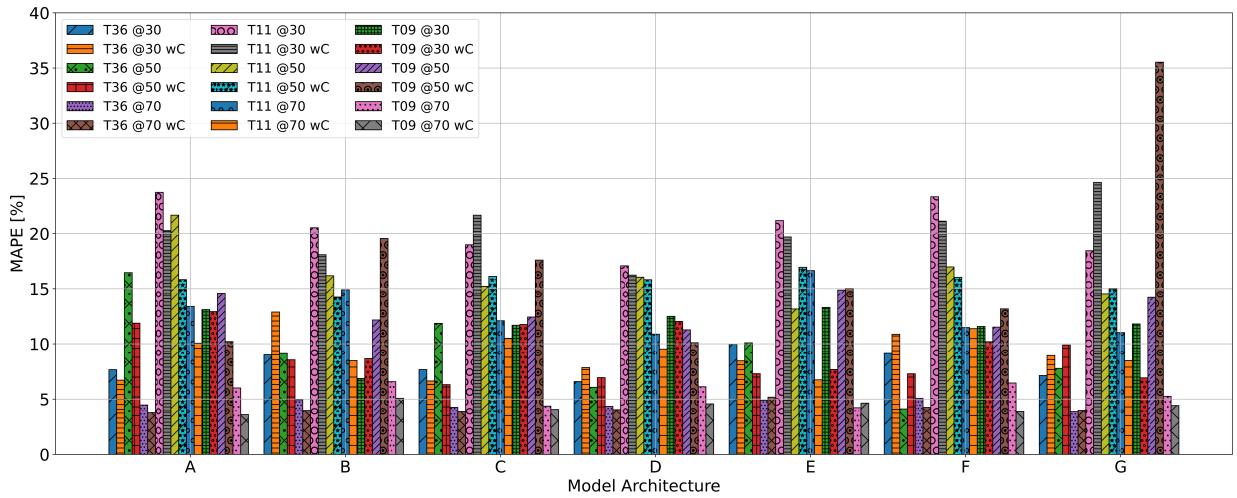


Figure 9: TFT architecture selection for long-term MOSFET ageing forecasting.

Figure 9 again confirms that the TFT hyperparameters impact on the obtained results. Among the architectures tested, the best TFT architecture was selected based on the best overall performance criteria. That is, it is selected by counting the number of times that each architecture performs best. It is observed that the best architecture is represented by model D.

5.2.2 Prediction Results

Finally, a LOO test was performed by training the model with 2 tests and leaving the third for the evaluation of the model. Accordingly, Figure 10 shows the obtained results for different (a) models: TFT, TFT with covariate (**TFT wCov**), and ELM; (b) LOO evaluation experiments: T36, T11, and T09; and (c) training proportions: 30%, 50%, 70% training samples, or equivalently, 70, 50 and 30 steps-ahead forecasting horizons, respectively.

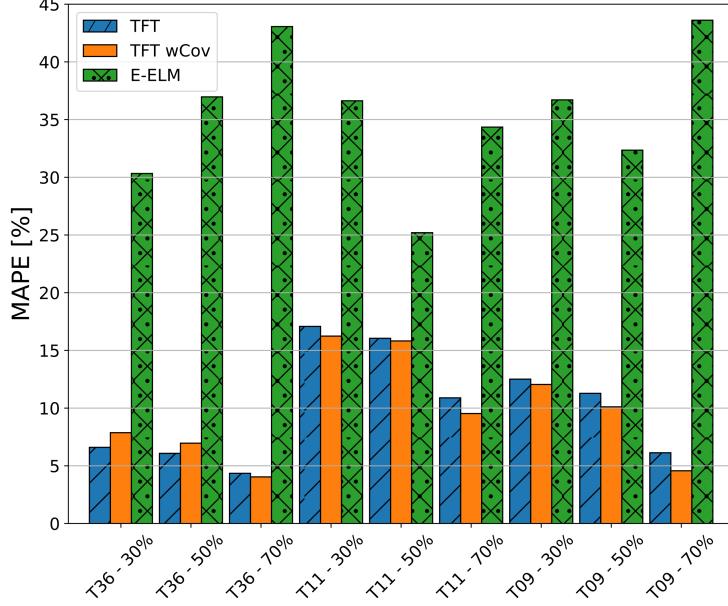


Figure 10: Long term $\Delta R_{DS_{ON}}$ prediction MAPE for different tests, models and training proportions.

Figure 10 shows that the TFT obtains a much smaller MAPE than the ELM for all the different configurations. Among the TFT variants, the TFT with covariate is the model with the lowest MAPE, except in two configurations (T36, 30% and 50%). As for the MAPE with respect to the prediction horizon, for the TFT architectures, the MAPE is reduced for smaller forecasting horizons.

Finally, taking the best TFT architecture Figure 11 shows prediction results with the 90% prediction intervals for the different tests and prediction horizons.

Figure 11 shows that the most accurate prediction results are obtained with the test 36 for all the tested prediction horizons. Furthermore, it can be observed that the prediction intervals are wider for tests 11 and 9, which is in agreement with the accuracy of the obtained results.

5.2.3 Attention mechanism

The transformer architecture enables the identification of attention points, which are indicators of the input sequences of the TFT architecture that are used to forecast (Vaswani et al., 2017). Focusing on 70-step ahead predictions (*i.e.*, 30% of training data), Figures 12 (a)-(c) show the attention curves for the different tests. Note that 70-step-ahead predictions are performed through concatenation. That is, 4 step-ahead predictions are used in an iterative manner to reach 70 step-ahead predictions.

Taking the maximum attention levels in Figures 12(a)-(c) with their corresponding time instant, Figures 12(d)-(f) show the identified attention points in the $R_{DS_{ON}}$ signal, which elucidate points where new failure mechanism arise or describe point of exacerbation. It can be observed that the TFT model pays attention to the main signal turning points for each test. In addition, it can be observed that, in general, as the prediction horizon progresses, the attention points become narrower and more focused.

6 Discussion

The final objective of the proposed comparative forecasting approach is to examine long- and short- term MOSFET ageing forecasting ability of classical state-space, statistical, and NN-based architectures, and novel transformer-architecture based forecasting models. This will enable to systematically select a prognostics model adapted to different forecasting needs. Yet, the selection of the final prognostics architecture is dependent on different factors and this section aims to discuss relevant criteria for an effective MOSFET ageing forecasting implementations.

6.1 Accuracy and computational cost

Short-term ageing forecasting results (cf. Table 3 and Figure 7) have shown that using classical state-space, statistical and NN-based architectures can be used to predict MOSFET ageing trajectory with E-ELM and ARIMA models obtaining the best results, with a good performance of the UKF model. The TFT model performs with an average accuracy for short-term forecasts, and therefore, for short-term forecasting the use of classical methods is suggested due to the simplicity with respect to the TFT model.

As for long-term ageing forecasts (cf. Figure 10), the TFT outperforms the best classical forecasting model (E-ELM). It is observed that the designed covariate informs the TFT model and improves prediction accuracy. Namely, the main strength of the TFT architecture is the capability to integrate additional covariates and learn influential dynamics that can be used to predict future ageing values. In contrast, classical state-space and statistical forecasting methods tend to accumulate errors for longer prediction horizons. Similarly, classical NN-based forecasting methods struggle to extrapolate beyond training data, and therefore increase the prediction error for long-term forecasting horizons. Therefore, TFT can be used as an offline approach to elucidate different failure mechanisms and evaluate long-term ageing trajectories.

The main differences in performance and computational cost are caused by the architectural complexity of the employed forecasting methods. Classical statistical forecasting methods, rely on autoregressive, moving average, and exponential smoothing techniques to track trends and correct errors. These methods are computationally efficient due to their parametric structure, but may struggle with capturing complex nonlinear degradation patterns.

State-space models require the specification of an underlying system model to track the data and update state estimates in real time. Their computational complexity increases with the dimensionality of the state-space and the nonlinearity of the system equations. Traditional data-driven forecasting methods rely on nonlinear transformations to approximate complex degradation trends. Their computational cost scales with the size of the network, the length of historical data sequences, and the training process.

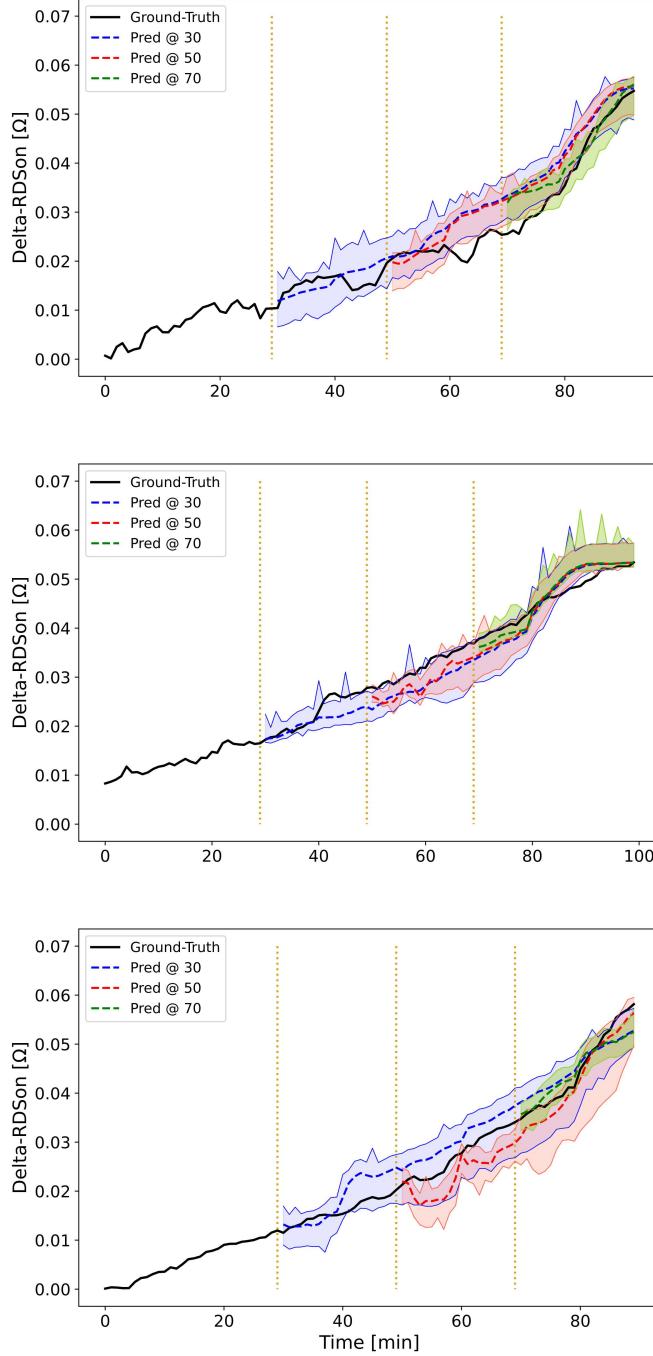


Figure 11: TFT forecasting results (mean and 90% prediction intervals) for tests 11 (top), 36 (middle) and 9 (bottom). Vertical dashed lines indicate forecasting instants using 30%, 50% and 70% of training data.

The TFT represents the most computationally demanding approach due to its self-attention mechanisms, gating layers, and multi-horizon forecasting architecture. Unlike other methods, the TFT can selectively focus on relevant past information through the attention mechanism, incorporate multi-source inputs through covariates, and provide uncertainty-aware long-term forecasts. The added complexity translates into higher computational cost, especially during training, but as observed in the obtained results, this leads to a superior accuracy when forecasting highly non-

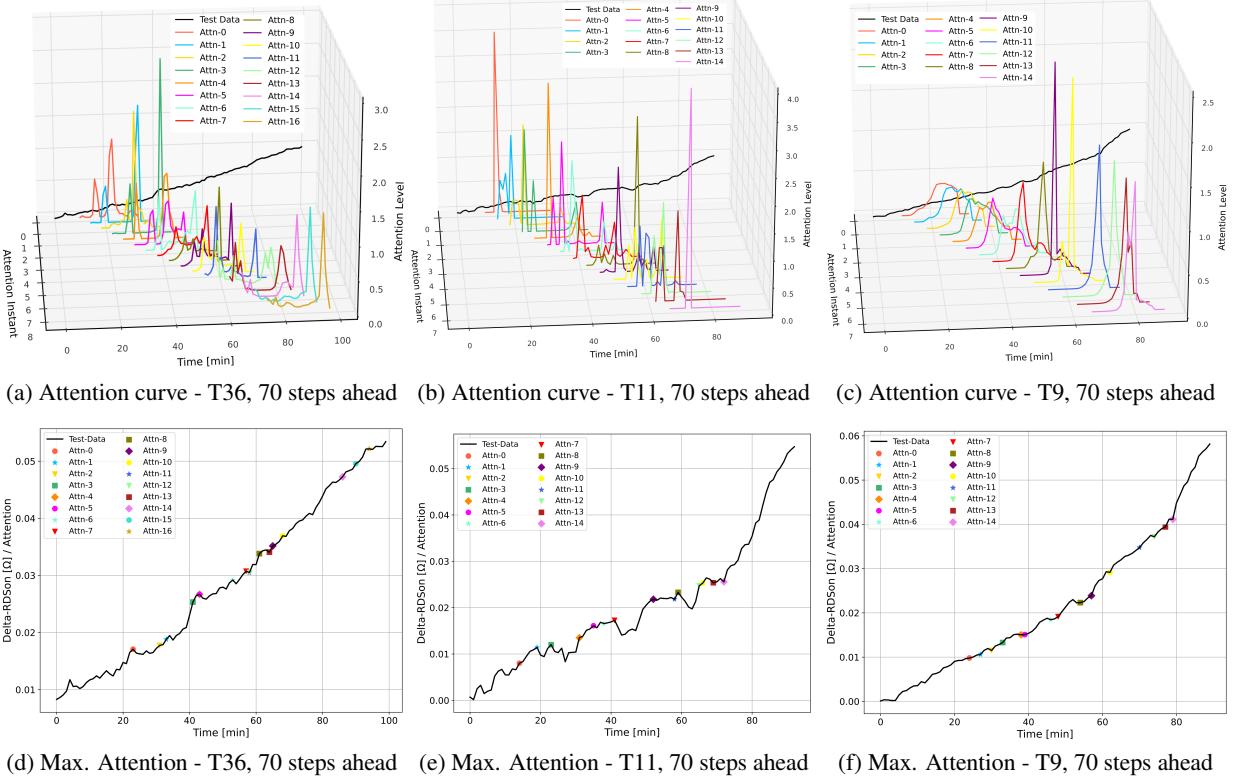


Figure 12: TFT attention plots for different tests. Top: Attention curves. Bottom: maximum attention points with respect to $R_{DS_{ON}}$.

linear and long-term degradation patterns. Table 6 displays a qualitative summary of the analysed methods, including their strengths, weaknesses and best application contexts.

Table 6: Qualitative comparison of the analysed methods.

Method	Strengths	Weaknesses	Best Use
ARIMA, Holt	Fast, simple, and interpretable. Good for linear trends. Low computational cost.	Difficulties with nonlinear trends	Short-term
EKF, UKF	Robust to nonlinearities. Real-time tracking	Sensitive to state-space. Moderate computational cost	Short-term
E-NN, E-ELM	Fast inference. Capture nonlinear and complex trends.	Requires large datasets. Moderate computational cost	Short-, Mid-term
TFT	Good long term accuracy. Handles long-term dependencies	Requires large datasets. Very high computational cost.	Long-term forecasting

6.2 TFT for Ageing Forecasting

The TFT architecture has the ability to incorporate covariates along with attention mechanisms that can be very useful for long-term prognostics. However, the adaptation of the TFT architecture for prognostics tasks is not trivial. This is mainly because relevant ad-hoc covariates need to be developed, which inform the TFT architecture about the evolution of the analysed ageing process.

6.3 MOSFET failure mechanisms

While the proposed approach has considered the exponential crack growth as the primary failure mechanism, there are other failure mechanisms present in the analysed ageing process, such as thermal runaway.

The impact of additional failure mechanisms becomes more significant beyond a certain key ageing threshold, making it challenging to track the ageing trajectory. This threshold can guide the selection of models. For well-defined modeling equations, such as exponential growth, where additional failure mechanisms emerge in later stages, statistical forecasting and tracking methods may suffice to forecast the ageing trajectory. Conversely, equations describing the failure physics that are ill-defined due to the interaction of multiple nonlinear failure mechanisms can benefit from data-driven models.

7 Conclusions

MOSFETs are key semiconductor-based power elements for various power-electronic systems. Their most critical failure mode is bond wire lift-off, which is caused by crack growth due to thermal fatigue. Focusing on this ageing mechanism, this research has evaluated the use of novel and classical forecasting methods for MOSFET ageing forecasting. Namely, Temporal Fusion Transformers (TFTs) have been adapted, implemented and tested for short- and long-term aging forecasting activities.

TFT models have been designed with ad-hoc covariates, and they have been compared with classical state-space tracking methods, including Extended and Unscented Kalman filters (EKF, UKF), classical forecasting methods, such as ARIMA and Holt, and other machine learning models used for forecasting, including ensemble models of Neural Networks (E-NN) and Extreme Learning Machines (ELM).

Obtained results show that, for short-term ageing forecasting horizons, the use of ARIMA and E-ELM models is the most appropriate. On the contrary, for long-term ageing prediction horizons, the use of TFT models greatly improves prediction accuracy with respect to all the tested models. **This highlights the capability of TFT models to incorporate (i) previous information for long-term future predictions making use of the attention mechanism and (ii) ageing-influencing data through covariates.**

The bond wire lift-off phenomena manifests as a exponential crack growth. The proposed aging forecasting approach allows obtaining an accurate modeling, which is very appropriate for short-term ageing prediction. With this model, all the analysed classic forecasting algorithms (EKF, UKF, ARIMA, Holt, E-ELM, E-NN) result in a proper short-term ageing estimation. **Given the relatively large number of computations of the TFT model, classic forecasting algorithms may be used for on-line prediction.**

For long-term prediction, classic forecasting algorithms are unable to obtain accurate results because, as the ageing of the MOSFET progresses, different failure modes interact and this is not captured. **In this context, TFTs represent an appropriate ageing forecasting alternative as they can incorporate multiple sources of information through covariates and focus on the most relevant information through the attention mechanism. To benefit from the TFT architecture for ageing forecasting, it is necessary to design and incorporate ageing covariates, which inform about the expected future inputs in the form of failure progression.**

The obtained results have shown that TFTs can obtain accurate forecasting results for long-term ageing forecasting. In addition, TFT attention points identify key ageing turning points, which are indicative of new failure modes or excessive accelerated ageing. All in all, it can be concluded that TFTs can be used to build MOSFET prognostics models. However, their computational cost may be an implementation limit, compared with classical algorithms, especially for on-line condition monitoring applications.

Future work could explore computationally efficient long-term ageing forecasting models, opening new opportunities in edge machine learning (Hua et al., 2023). This could enable low-latency, on-device inference in online monitoring systems for asset management (Garro et al., 2020). Another interesting direction is the adaptation of foundation models (Liang et al., 2024) for MOSFET ageing prediction. These large-scale pre-trained models could be fine-tuned for long-term degradation forecasting, potentially improving generalization across different operational conditions and component types.

Code Availability

The code will be available on the website: <https://github.com/joxeina/AgeingForecastingMOSFETs>.

Acknowledgements

This research was partially funded by the Basque Government, Department of Education (grant No. KK-2024/00030 and IT1504-22). J. I. Aizpurua is funded by the Ramón y Cajal Fellowship, Spanish State Research Agency (grant number RYC2022-037300-I), co-funded by MCIU/AEI/10.13039/501100011033 and FSE+.

References

- Li Zhu, Junghui Chen, and Chun-I Chen. Prognostics for semiconductor sustainability: Tool failure behavior prediction in fabrication processes. *IEEE Trans. Systems, Man, and Cybernetics: Systems*, pages 1–11, 2024. doi:10.1109/TSMC.2024.3359851.
- Ingeborg de Pater and Mihaela Mitici. Developing health indicators and rul prognostics for systems with few failure instances and varying operating conditions using a lstm autoencoder. *Engineering Applications of Artificial Intelligence*, 117:105582, 2023. ISSN 0952-1976. doi:<https://doi.org/10.1016/j.engappai.2022.105582>.
- Mengfan Zhang, Pere Izquierdo Gómez, Qianwen Xu, and Tomislav Dragicevic. Review of online learning for control and diagnostics of power converters and drives: Algorithms, implementations and applications. *Renewable and Sustainable Energy Reviews*, page 113627, 2023.
- M. Novak, A. Sangwongwanich, and F. Blaabjerg. Monte carlo-based reliability estimation methods for power devices in power electronics systems. *IEEE Open Journal of Power Electronics*, 2:523–534, 2021. doi:10.1109/OJPEL.2021.3116070.
- Abu Hanif, Yuechuan Yu, Douglas DeVoto, and Faisal Khan. A comprehensive review toward the state-of-the-art in failure and lifetime predictions of power electronic devices. *IEEE Trans. Power Electron.*, 34(5):4729–4746, 2019. doi:10.1109/TPEL.2018.2860587.
- Yongqiang Kang, Yue Pan, Luzhi Dang, Zhaoyun Wang, Yu Meng, and Shuaibing Li. Partial discharge characterization and mechanism of igt module power cycling aging processes. *IEEE Transactions on Components, Packaging and Manufacturing Technology*, pages 1–1, 2025. doi:10.1109/TCMPMT.2025.3531437.
- Ander Zubizarreta, Markel Penalba, David Garrido, Unai Markina, Xabier Ibarrola, and Jose Aizpurua. Uncertainty assessment framework for igt lifetime models. a case study of solder-free modules. *International Journal of Prognostics and Health Management*, 16(1), 2025.
- J. Lutz, H. Schlangenotto, U. Scheuermann, and R. De Doncker. *Semiconductor Power Devices: Physics, Characteristics, Reliability*. Springer International Publishing, 2018. ISBN 9783319709178. URL <https://books.google.co.uk/books?id=gaRMDwAAQBAJ>.
- Fraser Anderson, Karoline Pelka, Julia Walgern, Timo Lichtenstein, and Katharina Fischer. Trends and influencing factors in power-converter reliability of wind turbines: A deepened analysis. *IEEE Transactions on Power Electronics*, pages 1–12, 2025. doi:10.1109/TPEL.2025.3530163.
- José R. Celaya, Abhinav Saxena, Chetan S. Kulkarni, Sankalita Saha, and Kai Goebel. Prognostics approach for power mosfet under thermal-stress aging. In *Proc. Annual Reliability & Maintainability Symposium*, pages 1–6, 2012. doi:10.1109/RAMS.2012.6175487.
- Tomislav Dragičević, Patrick Wheeler, and Frede Blaabjerg. Artificial intelligence aided automated design for reliability of power electronic systems. *IEEE Trans. Power Electron.*, 34(8):7161–7171, 2019. doi:10.1109/TPEL.2018.2883947.
- Ning Wang, Yongbin Jiang, Weihao Hu, Yanbo Wang, and Zhe Chen. An ann-aided parameter design method for cllc-type dab converters considering parameter perturbation. *IEEE Transactions on Industrial Electronics*, pages 1–11, 2024. doi:10.1109/TIE.2024.3451135.
- Shuai Zhao, Frede Blaabjerg, and Huai Wang. An overview of artificial intelligence applications for power electronics. *IEEE Trans. Power Electron.*, 36(4):4633–4658, 2021. doi:10.1109/TPEL.2020.3024914.
- Carlos J. Alonso, Belarmino Pulido, Mario Cartón, and Anibal Bregón. A big data architecture for fault prognostics of electronic devices: Application to power mosfets. *IEEE Access*, 7:102160–102173, 2019. doi:10.1109/ACCESS.2019.2929111.
- Xiang Wang, Weiwei Wei, Yanhui Zhang, Wei Feng, Guoqing Xu, and An Xiang. A data-driven lifetime prediction method for thermal stress fatigue failure of power mosfets. *Energy Reports*, 8:467–473, 2022. ISSN 2352-4847. doi:10.1016/j.egyr.2022.10.137.

- Z. Li, Z. Zheng, and R. Outbib. A prognostic methodology for power mosfets under thermal stress using echo state network and particle filter. *Microelectronics Reliability*, 88-90:350–354, 2018. ISSN 0026-2714. doi:10.1016/j.microrel.2018.07.137.
- Mohammadreza Baharani, Mehrdad Biglarbegian, Babak Parkhideh, and Hamed Tabkhi. Real-time deep learning at the edge for scalable reliability modeling of si-mosfet power electronics converters. *IEEE Internet of Things Journal*, 6(5):7375–7385, 2019.
- Wenfa Kang, Sen Tan, Juan C Vasquez, Josep M Guerrero, Tobias Hertle, Thomas Gietzold, Andrew Benn, and Baoze Wei. A data-driven lifetime prediction method for thermally aged sic mosfet applications. In *2024 Prognostics and System Health Management Conference (PHM)*, pages 281–286. IEEE, 2024.
- Civan Lezgin Kahraman, Darius Roman, Lucas Kirschbaum, David Flynn, and Jonathan Swingler. Machine learning pipeline for power electronics state of health assessment and remaining useful life prediction. *IEEE Access*, 2024.
- Qian Yang, Muhammed Ali Gultekin, Vahe Seferian, Krishna Pattipati, Ali M. Bazzi, Francesco A. N. Palmieri, Ravi Rajamani, Shailesh N. Joshi, Muhammed Farooq, and Hiroshi Ukegawa. Incipient residual-based anomaly detection in power electronic devices. *IEEE Trans. Power Electronics*, 37(6):7315–7332, 2022. doi:10.1109/TPEL.2022.3140721.
- Shuai Zhao, Viliam Makis, Shaowei Chen, and Yong Li. Health assessment method for electronic components subject to condition monitoring and hard failure. *IEEE Trans. Instrum. Meas.*, 68(1):138–150, 2018.
- Qunfang Wu, Boyuan Xu, Lan Xiao, and Qin Wang. A remaining useful life prediction method of sic mosfet considering failure threshold uncertainty. *IET Power Electronics*, 17(12):1594–1606, 2024.
- Alireza Alghassi, Suresh Perinpanayagam, and Mohammad Samie. Stochastic rul calculation enhanced with tdnn-based igit failure modeling. *IEEE Transactions on Reliability*, 65(2):558–573, 2016. doi:10.1109/TR.2015.2499960.
- Marco Rigamonti, Piero Baraldi, Allegra Alessi, Enrico Zio, Daniel Astigarraga, and Ainhoa Galarza. An ensemble of component-based and population-based self-organizing maps for the identification of the degradation state of insulated-gate bipolar transistors. *IEEE Transactions on Reliability*, 67(3):1304–1313, 2018. doi:10.1109/TR.2018.2834828.
- Dengyu Xiao, Chengjin Qin, Jianwen Ge, Pengcheng Xia, Yixiang Huang, and Chengliang Liu. Self-attention-based adaptive remaining useful life prediction for igit with monte carlo dropout. *Knowledge-Based Systems*, 239:107902, 2022. ISSN 0950-7051. doi:10.1016/j.knosys.2021.107902.
- Shuhan Deng, Zhuyun Chen, Hao Lan, Ke Yue, Zhicong Huang, and Weihua Li. Remaining useful life prediction with spatio-temporal graph transform and weakly supervised adversarial network: An application in power components. *Energy*, 313:133599, 2024. ISSN 0360-5442. doi:<https://doi.org/10.1016/j.energy.2024.133599>.
- Zhonghai Lu, Chao Guo, Mingrui Liu, and Rui Shi. Remaining useful lifetime estimation for discrete power electronic devices using physics-informed neural network. *Scientific Reports*, 13(1):10167, 2023.
- Youssef Fassi, Vincent Heiries, Jerome Boutet, and Sebastien Bousseau. Toward physics-informed machine-learning-based predictive maintenance for power converters—a review. *IEEE Transactions on Power Electronics*, 39(2):2692–2720, 2024. doi:10.1109/TPEL.2023.3328438.
- Konstantinos Benidis, Syama Sundar Rangapuram, Valentin Flunkert, Yuyang Wang, Danielle Maddix, Caner Turkmen, Jan Gasthaus, Michael Bohlke-Schneider, David Salinas, Lorenzo Stella, François-Xavier Aubet, Laurent Callot, and Tim Januschowski. Deep learning for time series forecasting: Tutorial and literature survey. *ACM Comput. Surv.*, 55(6), dec 2022. ISSN 0360-0300. doi:10.1145/3533382.
- Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. *Philosophical Trans. of the Royal Society A: Mathematical, Physical & Eng. Sciences*, 379(2194):20200209, 2021. doi:10.1098/rsta.2020.0209.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- William Gilpin. Large-scale statistical forecasting models reassess the unpredictability of chaotic systems, 2023.
- Luca Biggio, Tommaso Bendinelli, Chetan Kulkarni, and Olga Fink. Ageing-aware battery discharge prediction with deep learning. *Applied Energy*, 346:121229, 2023. ISSN 0306-2619. doi:10.1016/j.apenergy.2023.121229.
- V.C. Valchev and A. Van den Bossche. *Inductors and Transformers for Power Electronics (1st ed.)*. Springer International Publishing, 2005. ISBN 9781315221014. URL <https://doi.org/10.1201/9781420027280>.

- Serkan Dusmez, Hamit Duran, and Bilal Akin. Remaining useful lifetime estimation for thermally stressed power mosfets based on on-state resistance variation. *IEEE Trans. Ind. Appl.*, 52(3):2554–2563, 2016. doi:10.1109/TIA.2016.2518127.
- Charles K Chui and Guanrong Chen. *Kalman Filtering with real-time applications*. Springer, 2017.
- Eric A. Wan and Rudolph van der Merwe. *The Unscented Kalman Filter*, chapter 7, pages 221–280. John Wiley & Sons, Ltd, 2001. ISBN 9780471221548. doi:<https://doi.org/10.1002/0471221546.ch7>.
- Robin John Hyndman and George Athanasopoulos. *Forecasting: Principles and Practice*. OTexts, Australia, 2nd edition, 2018.
- Rob Hyndman, George Athanasopoulos, Christoph Bergmeir, Gabriel Caceres, Leanne Chhay, Mitchell O’Hara-Wild, Fotios Petropoulos, Slava Razbash, Earo Wang, and Farah Yasmeen. *forecast: Forecasting functions for time series and linear models*, 2024.
- Simone Scardapane. *Alice’s Adventures in a Differentiable Wonderland – Volume I*. 2024. URL <https://arxiv.org/abs/2404.17625>.
- Nikolaos Kourentzes, Devon K. Barrow, and Sven F. Crone. Neural network ensemble operators for time series forecasting. *Expert Systems with Applications*, 41(9):4235–4244, 2014. ISSN 0957-4174. doi:10.1016/j.eswa.2013.12.011.
- Xiaoqian Wang, Rob J Hyndman, Feng Li, and Yanfei Kang. Forecast combinations: An over 50-year review. *Int. J. Forecasting*, 39(4):1518–1547, 2023.
- Nikolaos Kourentzes. *nnfor: Time Series Forecasting with Neural Networks*, 2023. R package version 0.9.9.
- Bryan Lim, Sercan O. Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *Int. J. Forecasting*, 37(4):1748–1764, 2021. ISSN 0169-2070. doi:10.1016/j.ijforecast.2021.03.012.
- Julien Herzen, Francesco LÄ¶ssig, Samuele Giuliano Piazzetta, Thomas Neuer, LÄ©o Tafti, Guillaume Raille, Tomas Van Pottelbergh, Marek Pasieka, Andrzej Skrodzki, Nicolas Huguenin, Maxime Dumonal, Jan KoÅ»cisz, Dennis Bader, FrÃ©dÃ©rick Gusset, Mounir Benhedi, Camila Williamson, Michal Kosinski, Matej Petrik, and GaÃ«l Grosch. Darts: User-friendly modern machine learning for time series. *J. Machine Learning Research*, 23(124):1–6, 2022.
- Sariel Har-Peled and Benjamin Raichel. The fréchet distance revisited and extended. *ACM Transactions on Algorithms (TALG)*, 10(1):1–22, 2014.
- Haochen Hua, Yutong Li, Tonghe Wang, Nanqing Dong, Wei Li, and Junwei Cao. Edge computing with artificial intelligence: A machine learning perspective. *ACM Computing Surveys*, 55(9):1–35, 2023.
- Unai Garro, Eñaut Muxika, Jose Ignacio Aizpurua, and Mikel Mendicute. Fpga-based stochastic activity networks for online reliability monitoring. *IEEE Transactions on Industrial Electronics*, 67(6):5000–5011, 2020. doi:10.1109/TIE.2019.2928244.
- Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’24, page 6555–6565. Association for Computing Machinery, 2024. ISBN 9798400704901. doi:10.1145/3637528.3671451.