

PowerBrain: An automatic data extraction tool for semiconductor datasheets

Fanghao Tian ¹, Qingcheng Sui ¹, Jiaze Kong ¹, Wilmar Martinez ¹

¹ EnergyVille - KU Leuven, Belgium

Corresponding author: Fanghao Tian, fanghao.tian@kuleuven.be

Abstract

Design automation (DA) for power converters has become an emerging field of research due to the design complexity. One of the primary challenges is modelling and comparing the power losses for an enormous number of components. This paper introduces a new and comprehensive artificial intelligence (AI) tool called PowerBrain, which is available online and is regularly updated and improved through the website <https://www.powerbrain.ai>. The tool is designed to extract nonlinear dynamic features from semiconductor device datasheets. The results indicate that the tool, specifically designed for power transistor device datasheets, can precisely extract the data within 2% relative error, resulting in a significant decrease in the time required for transistor data collection. Ultimately, it facilitates the creation of a standardized database, aiding power converter designers in the development processes.

1 Introduction

The increasing need for power electronic converters, driven by the adoption of electric mobility and renewable energy industries, requires power converters that adhere to strict standards of efficiency, density, and durability. The wide range of metal oxide semiconductor field effect transistors (MOSFET) and the introduction of advanced semiconductors such as silicon carbide (SiC) and gallium nitride (GaN) make it challenging to choose the most suitable power switches due to their various performance characteristics in terms of power loss, size, and cost [1]. Accurate preselection relies on computational assessments of power loss in silicon, SiC MOSFETs, and GaN HEMTs. This is because power loss in these materials is dynamic and non-linear in nature. Recent research has concentrated on analytical loss models that use dynamic datasheet data, highlighting the importance of variables such as the relationship between transconductance and channel current, as well as the impact of threshold voltage hysteresis (TVH) [2] - [5]. These models, which have been confirmed by experimental data with acceptable margins of error, indicate that a comprehensive dynamic data database could greatly assist in the efficient assessment and choice of power switches among the wide range of possibilities available in the market. In this case, an exclusive database which contains all the characteristics of thousands of components

needs to be constructed.

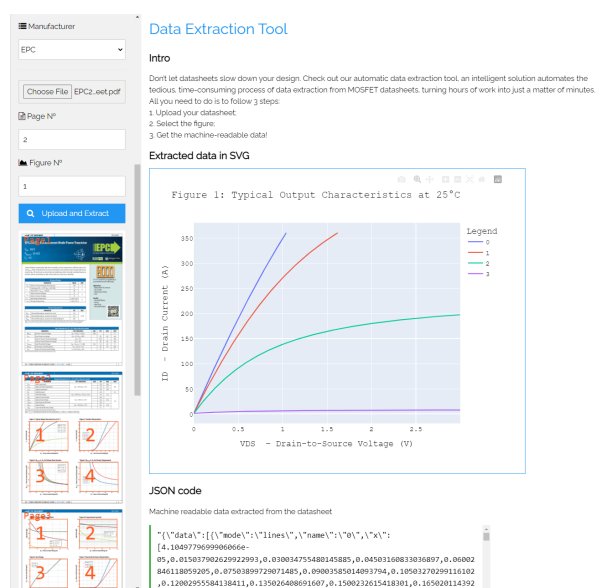


Fig. 1: A screenshot of the PowerBrain data extraction tool.

However, it is impractical to manually extracting the dynamic data from datasheets since it is inefficient and prone to inaccuracies. However, the integration of artificial intelligence (AI) within power electronics has facilitated the automation of design processes. Specifically, convolutional neural network (CNN) algorithms excel in image processing, facilitating the process of extracting dynamic data from datasheets. This necessitates the development of a specialized tool tailored for high-precision

data extraction from line charts in datasheets. Consequently, this paper introduces an AI-powered tool aimed at automating data extraction from MOSFET and GaN HEMT datasheets, thus supporting the creation of a database to improve power loss modeling for power converter design. Designed as a comprehensive solution, PowerBrain streamlines the extraction process, allowing designers to input a PDF datasheet and automatically retrieve necessary dynamic data, as shown in Fig. 1. In the end, a demonstration of using datasheet data to calculate the power losses are also introduced with a case of comparing multiple devices.

2 PowerBrain

PowerBrain is a streamlined tool that can extract dynamic data directly from a line chart by inputting the original datasheet file in PDF format. The process is divided into 2 major parts. Before analyzing individual figures, they are recognized and separated from the source PDF pages. Next, the dynamic data contained in individual figures are extracted into a machine-readable format.

Once a PDF datasheet is uploaded, a preview will display all of the figures contained in the file. Next, the individual figure can be analyzed with the data extracted into a machine readable format.

2.1 CNN based key elements detection

Throughout the workflow of PowerBrain, CenterNet [6], an algorithm for detecting objects, is used to locate figures in PDF pages and identify the key elements inside those figures, which helps with extracting line data in the following steps. Thus, 2 CenterNet models are trained.

Manufacturers offer a variety of formats for their datasheets. Certain pages have six figures, whilst others only have four. In order to enhance the universality of the object detection algorithm, we gather datasheets from multiple manufacturers with diverse layouts and manually annotate them. The dataset consists of 2000 screenshots of datasheet pages, including annotations for all line charts. Additionally, there is a separate dataset consisting of 3800 individual figures that have been extracted from datasheets. These figures contain annotations for six important elements: title, legend, coordinate origin corner, label, x-axis and y-axis tick numbers, and other relevant information. Fig.2 displays some examples of annotated figures.

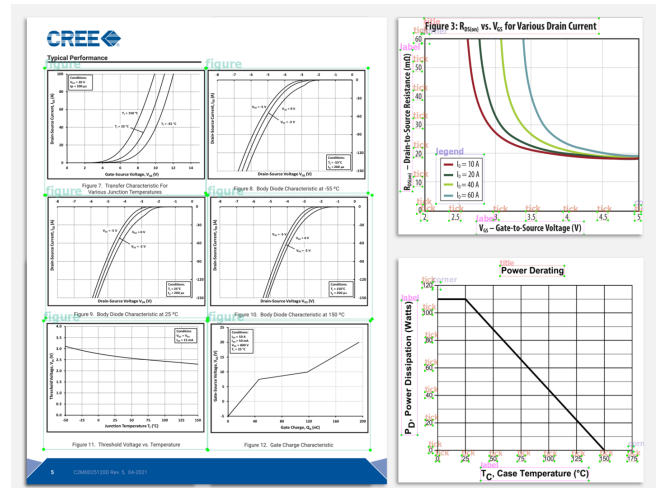


Fig. 2: Examples of the annotated PDF page and figures for object detection algorithm.

The CenterNet models are trained after constructing the dataset. Firstly, the dataset is randomly partitioned into training, validation, and test subsets, comprising 80%, 10%, and 10% of the total data, respectively. This division ensures a robust evaluation of the model's performance across different stages of training and validation. Then 300 epochs of training are conducted. In order to accelerate the training process, the models incorporate ResNet50 model [7], utilizing parameters derived from a large benchmark dataset, as the backbone CNN. In the first 150 epochs, the ResNet50 parameters remain fixed, providing a stable foundation for training. Subsequently, the parameters become adjustable for the final 150 epochs, allowing for fine-tuning and enhanced learning. This approach leverages the pre-trained ResNet50's capabilities, thereby reducing the necessity for extensive training on the backbone CNN.

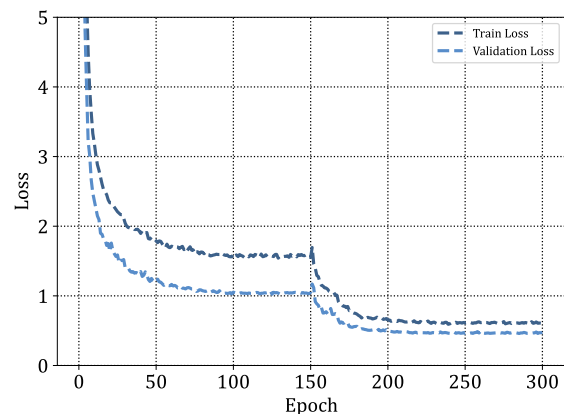


Fig. 3: Convergence curve for model 1 training.

After 300 epochs of training, the training losses

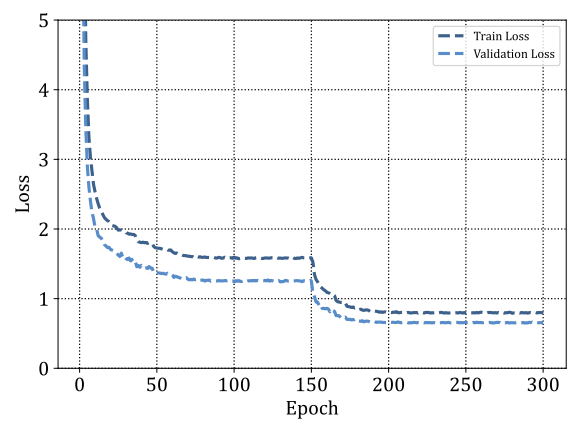


Fig. 4: Convergence curve for model 2 training.

and validation losses for model 1 reach 0.6064 and 0.4746, respectively. Meanwhile, for model 2, the losses reach 0.7909 and 0.6533. Fig.3 and Fig.4 show the convergence curves for model 1 and model 2.

As is shown in Fig.3 and Fig.4, the training and validating losses on both models decrease gradually in the first 150 epochs with the pre-trained backbone CNN being frozen. However, the losses reach a steady state at around 1-2. After the backbone CNN being unfreezed, the losses decreased further and reach a value smaller than 1. It shows that the training process was accelerated by the frozen scheme. In addition, given the small losses, the two models can both detect the target objects in the image well. Fig.5 shows some examples of the detection results on both models.

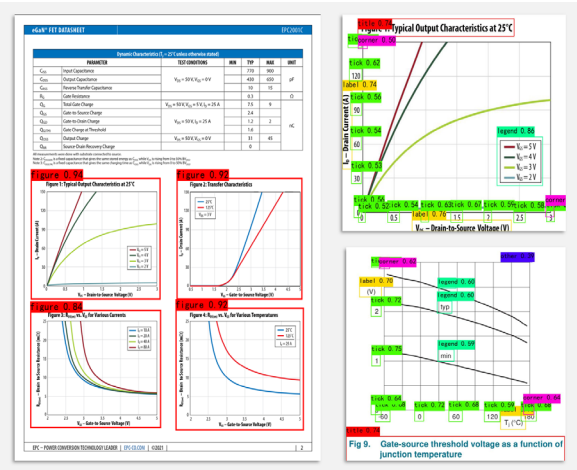


Fig. 5: Examples of the object detection results on PDF page and figures.

Finally, in order to evaluate the overall performance of the models, the test dataset was used to manually check the performance. For model 1, if all of the figures on the page is detected, the result is

regarded as accurate. For model 2, if all of the key elements are detected, the result is correct. Based on this rule, the results are listed in TABLE 1.

Tab. 1: Test results of two CenterNet models

	test data size	accuracy rate
Model one	200	97.5%
Model two	380	94.5 %

The results show that both of the trained models for detecting the figures from PDF page and detecting key elements in a figure can detect the objects in accurately. In the next section, the detected information will be uzlized for further process on extracting the line curve.

Nonetheless, there are still some failure cases. For model 1, some figures were missing. In the mean-time, the failure cases are more diverse in model 2. It includeds the cases when short labels are misidentified as tick numbers, an extra corner detected on the grid, and legends are mixed with other information. Luckily, not all of them compromise the following data extraction process. Since multiple verifications are conducted in the following steps, such as an additional filter on tick numbers, relationship checking among the tick numbers, and etc.

2.2 Line data extraction

PowerBrain is a single-flow tool for extracting the dynamic data from the datasheets in PDF format. After CenterNet Model 1 detects the figures from original PDF pages, the individual figures are cut out and passed to CenterNet model 2, through which the six key elements are detected, which are title, legend, label, tick numbers, corner, and other additional information.

Based on the detection result from model 2, three steps follow to extract the line data. Firstly, the texts are recognized by the Optical Character Recognition (OCR) engine Tesseract [9], including not only the normal characters in title, label, and legends, but also the numbers. Secondly, traditional image processing algorithms are applied for extracting the lines in pixel coordination. Lastly, the background grids are detected and matched with tick numbers to convert the result from pixels coordination to number coordination.

Fig.6 depicts the results of the extracted line in the width of 1 pixel under the coordination of pixels and the final results in number coordination. More

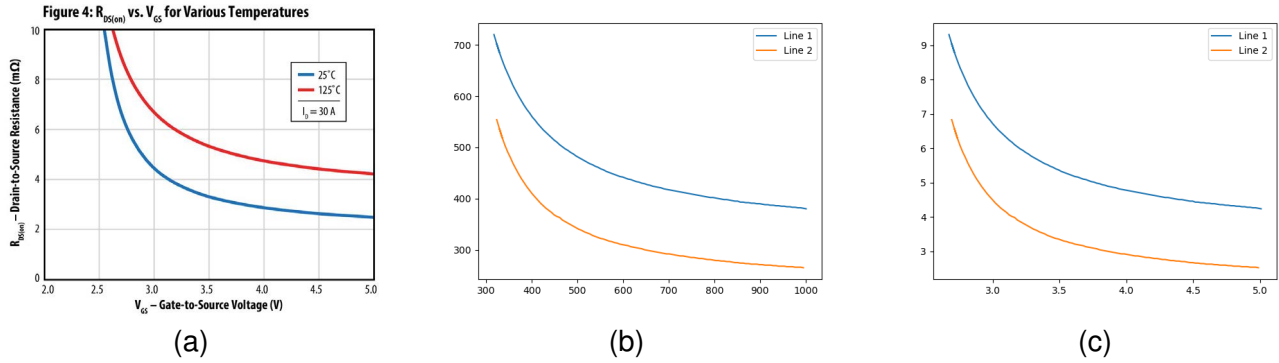


Fig. 6: Line extraction process (a) Original image. (b) Result in pixel coordination. (c) Result in number coordination.

details can be found in [8]. Afterwards, 100 figures were tested and compared with the real data. The results show an average relative error of 1.6%, illustrating the accuracy of the extracted data.

3 A case study: using datasheet data collected by PowerBrain for preselection

In the process of designing power electronics, the selection of semiconductor components has a substantial effect on the performance of power converters. Generally, engineers choose components based on their experience. It is impractical to manually analyze tens of thousands of semiconductor components, and the market is filled with semiconductor switching devices with comparable voltage and current characteristics. Consequently, more time is spent selecting a component.

Moreover, the factors and variables of semiconductors are of significant relevance. A precise and efficient model is necessary to identify which of them or which combination of them can increase the overall efficiency of the converter. In this section, a case study of using PowerBrain data for selecting the components is conducted.

One of the major criteria of choosing a MOSFET or GaN HEMT component is the power losses under the specific application conditions. However, the power losses are dependent on many parameters that are varying under different working conditions such as drain current, temperature, and etc.

In the case of conduction losses, they are related to the drain current $I_{DS(rms)}$ and the on-resistance $R_{DS(on)}$.

$$P_{cond} = I_{DS(rms)}^2 R_{DS(on)} \quad (1)$$

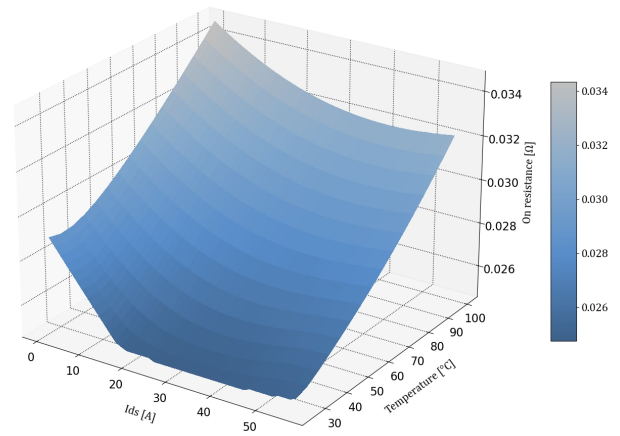


Fig. 7: An interpolation map of on-resistance.

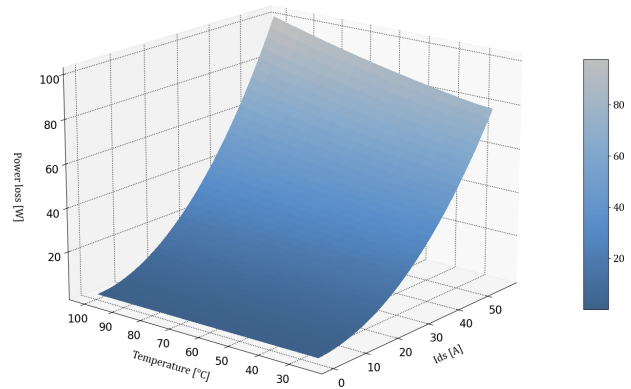


Fig. 8: An look-up table on conduction losses.

However, the $R_{DS(on)}$ is highly sensitive to operation conditions, such as gate-to-source voltage, drain current, and temperature. All of the specific dynamic information are contained in the figure of On-resistance vs. Drain Current (For various temperatures), On-Resistance vs. Temperature (For Various Gate Voltage), and On-resistance vs. Temperature.

All those figures can be easily extracted from datasheet by PowerBrain, and then an interpola-

tion is applied to obtain an accurate function of On-resistance under various operating conditions. If we suppose the gate voltage is determined, the interpolation result is shown in Fig.7. Furthermore, an conduction losses map can also be constructed given the varying drain current and temperatures, as is shown in Fig.8.

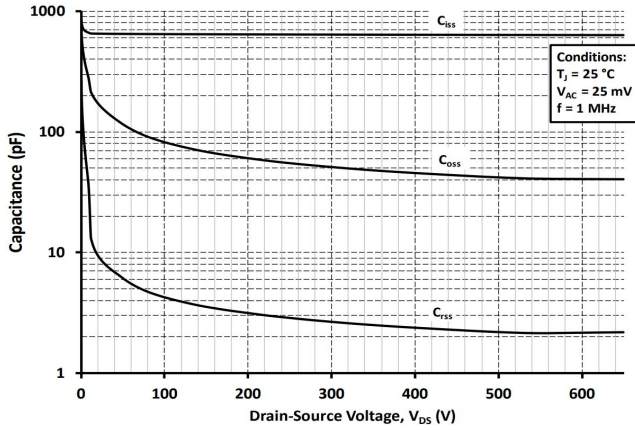


Fig. 9: An example of parasitic capacitance vs. voltage.

Given that parasitic capacitances vary with voltage, the equivalent capacitance must be obtained by integrating the parasitic capacitance curve from the datasheet, as demonstrated in (2).

$$C_{x,eq} = \frac{1}{V_o} \int_0^{V_o} C_x(v) dv, x = iss, oss, rss \quad (2)$$

Another example is the transconductance g_m , which is a function of the channel current as (3).

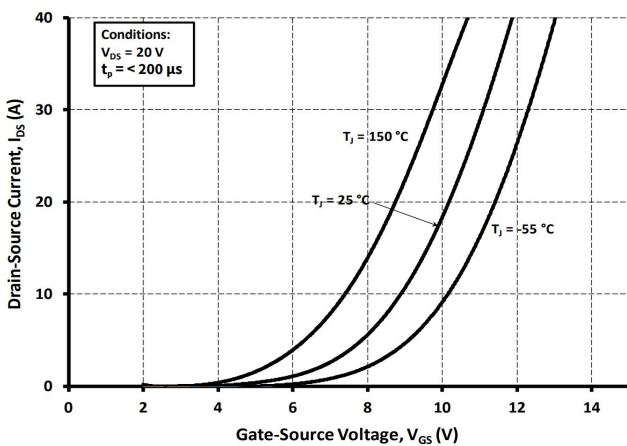


Fig. 10: An example of transfer characteristic.

$$i_{ch} = k_1 (v_{gs} - V_{th})^x + k_2 \quad (3)$$

In which the constants k_1 , k_2 , and x are obtained by curve fitting on Fig. 10.

In summary, all the dynamic data that is needed for calculating the power losses can be extracted from datasheet easily by PowerBrain. It allows for the evaluation of numerous components under specific application conditions without the need for experimental testing and individual comparisons. The complete power loss model by using datasheet data can also be found in [3] and [10].

Finally, a demonstration of rapid comparisons among multiple devices using data collected from datasheets is conducted. Consider an operating condition of $V_{ds} = 50V$ and $I_{ds} = 5A$ in a half-bridge configuration with a duty cycle of 0.5 and a frequency of 50 kHz. A comparison of power losses across 30 devices is illustrated in Fig.11.

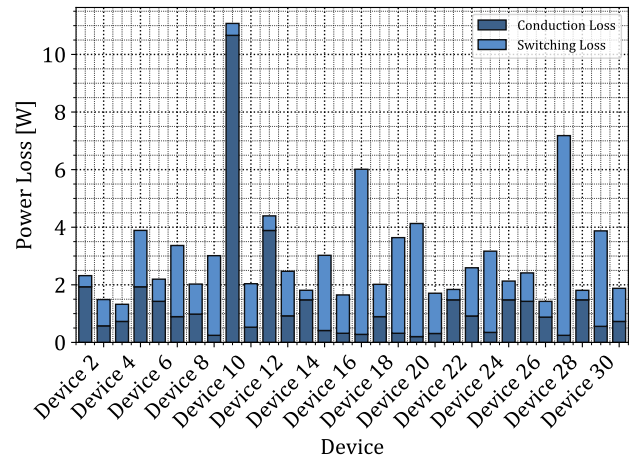


Fig. 11: A comparison of power losses on 30 components.

It is clear that some devices have a large on-resistance which produce high conduction losses while some devices generate more switching losses. In summary, a database compiled from the datasheets of various semiconductor devices can significantly enhance the evaluation and comparison of multiple MOSFETs. This resource aids power electronics designers in selecting the most suitable device based on their specific design requirements.

4 Conclusion

Power electronics designers need datasheets for accurate and current semiconductor component information. However, obtaining data from datasheet graphics is challenging, particularly when comparing numerous similar components. In addition, the data accuracy is essential for power loss model

analysis. And the component library is rapidly expanding as the semiconductor industry grows daily. Thus, a dedicated automatic data extraction tool for semiconductor devices, such as MOSFETs, GaN HEMTs, IGBT, and others, is essential to automate the work of power electronics designers.

PowerBrain offers an automated data extraction tool powered by AI algorithms. The dynamic data extracted by this tool can enhance the accuracy of power loss calculation models, thereby improving the optimization of power converter designs. This tool facilitates the rapid construction of a comprehensive dynamic database for the design and automation of power electronics. Based on the simulation results, powerbrain can automatically extract the data directly from the PDF file within 20s per figure, and the relative error of the data is less than 2%.

However, variations in datasheets from different manufacturers pose a challenge in creating a universal, reliable, and accurate database. Legends and lines are hard to match because manufacturers use different styles. PowerBrain addresses the above difficulties with customized algorithms for different manufacturers. However, standardizing datasheets across industries is crucial for introducing AI into power electronics designs and improving design automation for the future.

References

- [1] Mouser Electronics. Accessed on: March 18th, 2024. [Online]. Available: <https://eu.mouser.com/c/semiconductors/discrete-semiconductors>.
- [2] X. Huang, Q. Li, Z. Liu and F. C. Lee, "Analytical Loss Model of High Voltage GaN HEMT in Cascode Configuration," in *IEEE Transactions on Power Electronics*, vol. 29, no. 5, pp. 2208-2219, May 2014, doi: 10.1109/TPEL.2013.2267804.
- [3] D. Christen and J. Biela, "D. Christen and J. Biela, "Analytical Switching Loss Modeling Based on Datasheet Parameters for mosfets in a Half-Bridge," in *IEEE Transactions on Power Electronics*, vol. 34, no. 4, pp. 3700-3710, April 2019, doi: 10.1109/TPEL.2018.2851068.
- [4] C. Qian, Z. Wang, G. Xin and X. Shi, "Datasheet Driven Switching Loss, Turn-ON/OFF Overvoltage, di/dt, and dv/dt Prediction Method for SiC MOSFET," in *IEEE Transactions on Power Electronics*, vol. 37, no. 8, pp. 9551-9570, Aug. 2022, doi: 10.1109/TPEL.2022.3152529.
- [5] N. Wang, J. Zhang and F. Deng, "Improved SiC MOSFET Model Considering Channel Dynamics of Transfer Characteristics," in *IEEE Transactions on Power Electronics*, vol. 38, no. 1, pp. 460-471, Jan. 2023, doi: 10.1109/TPEL.2022.3200456.
- [6] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang and Q. Tian, "CenterNet: Keypoint Triplets for Object Detection," 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea (South), 2019, pp. 6568-6577, doi: 10.1109/ICCV.2019.00667.
- [7] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770-778, doi: 10.1109/CVPR.2016.90.
- [8] F. Tian, D. B. Cobaleda and W. Martinez, "Automatic Data Extraction Based on Semiconductor Datasheet for Design Automation of Power Converters," 2022 International Power Electronics Conference (IPEC-Himeji 2022- ECCE Asia), Himeji, Japan, 2022, pp. 922-927, doi: 10.23919/IPEC-Himeji2022-ECCE53331.2022.9806859.
- [9] R. Smith, "An Overview of the Tesseract OCR Engine," Ninth International Conference on Document Analysis and Recognition (ICDAR 2007), Curitiba, Brazil, 2007, pp. 629-633, doi: 10.1109/ICDAR.2007.4376991.
- [10] F. Tian, S. Li, X. Ning, D. B. Cobaleda and W. Martinez, "Embedding-Encoded Artificial Neural Network Model for MOSFET Preselection: Integrating Analytic Loss Models with Dynamic Characteristics from Datasheets," 2024 IEEE Applied Power Electronics Conference and Exposition (APEC), Long Beach, CA, USA, 2024, pp. 1574-1580, doi: 10.1109/APEC48139.2024.10509354.