

Data-Driven Models for Yacht Hull Resistance Optimization

Exploring Geometric Parameters Beyond the Boundaries of the Delft Systematic Yacht Hull Series

Walker, Jake M.; Coraddu, Andrea; Oneto, Luca

DOI

[10.1109/ACCESS.2024.3404495](https://doi.org/10.1109/ACCESS.2024.3404495)

Publication date

2024

Document Version

Final published version

Published in

IEEE Access

Citation (APA)

Walker, J. M., Coraddu, A., & Oneto, L. (2024). Data-Driven Models for Yacht Hull Resistance Optimization: Exploring Geometric Parameters Beyond the Boundaries of the Delft Systematic Yacht Hull Series. *IEEE Access*, 12, 76102-76120. Article 10537207. <https://doi.org/10.1109/ACCESS.2024.3404495>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Received 9 April 2024, accepted 12 May 2024, date of publication 23 May 2024, date of current version 4 June 2024.

Digital Object Identifier 10.1109/ACCESS.2024.3404495

RESEARCH ARTICLE

Data-Driven Models for Yacht Hull Resistance Optimization: Exploring Geometric Parameters Beyond the Boundaries of the Delft Systematic Yacht Hull Series

JAKE M. WALKER^{ID}¹, (Member, IEEE), ANDREA CORADDU^{ID}¹, (Member, IEEE),
AND LUCA ONETO^{ID}², (Senior Member, IEEE)

¹Faculty of Mechanical Engineering, Delft University of Technology, 2628 CD Delft, The Netherlands

²Department of Informatics, Bioengineering, Robotics and Systems Engineering, Università degli Studi di Genova, 16126 Genova, Italy

Corresponding author: Jake M. Walker (j.m.walker@tudelft.nl)

ABSTRACT Optimizing vessel hull resistance is pivotal for enhancing maritime performance and minimizing environmental impacts. Traditional methods combine expert intuition with Data-Driven Models (DDMs), relying on parametrization to predict and optimize hull geometries using Experimental Fluid Dynamics (EFD) or Computational Fluid Dynamics (CFD) data. However, these conventional approaches are hampered by several limitations: they require significant human input, are computationally intensive and costly, and lack flexibility in adapting to new families of geometries or parameters beyond predefined ranges. Addressing these challenges, our research introduces a novel method that significantly reduces the need for human intervention, computational resources, and costs, while also improving the model's adaptability. By proposing a new parametrization technique that accurately encompasses the Delft Systematic Yacht Hull Series (DSYHS), we demonstrate that DDMs can be effectively trained directly on EFD datasets. This eliminates the dependency on extensive CFD simulations or the generation of new EFD data tailored to a specific investigation. Our approach matches the performance of leading-edge CFD models, even in extrapolating conditions, with physical plausibility and minimal human oversight. The validation of our method under various and increasingly complex extrapolating scenarios, employing statistical analyses on the DSYHS EFD dataset and comparisons with state-of-the-art CFD models, underscores the effectiveness of our proposal. Furthermore, we demonstrated that our model can successfully optimize hull resistance when navigating geometric parameters outside the confines of the DSYHS validating our results through leading-edge CFD simulations. This work addresses the limitations of existing methodologies by offering a novel approach more accurate, efficient, cost-effective, flexible, automated, and robust to extrapolation for hull resistance optimization.

INDEX TERMS Computational fluid dynamics, data-driven models, DelftBlue, Delft Systematic Yacht Hull Series, extrapolation, hull parametrization, hull resistance, optimization, sailing yachts.

The associate editor coordinating the review of this manuscript and approving it for publication was Utku Kose^{ID}.

I. INTRODUCTION

Vessel hull resistance optimization is a critical design problem [1], [2], [3]. The hull-form resistance must be minimized

to improve performance and, in the case of motorized vessels, to minimize the environmental footprint [4], [5]. In order to achieve this goal, a wide and multidimensional design space needs to be explored [6], which is not a trivial task given the complexity of the representation space [7], [8].

Conventional methods to assess the performance of candidate vessel designs are computationally intensive or time-consuming or both [8], [9], and [10]. In fact, the classic approach to determine the hull resistance is to perform Experimental Fluid Dynamics (EFD) using model scale tests [11], [12], [13]. However, considerable effort is required to construct a model scale of the candidate hull and to perform the test in the appropriate facility. For this reason, modern approaches rely on virtual experiments using Computational Fluid Dynamics (CFD) [1], [14], [15], [16], [17], [18], [19], [20], [21], [22]. CFD usually provides accurate results that can be validated via EFD to improve the trustworthiness of the virtual experiment [1], [17], [23]. Nevertheless, when it comes to optimizing the design of the vessel hull, assessing the performance of many different candidate designs is required [24]. In this setting, using CFD results is impractical due to its computational requirements [8], [9], [25]. For this reason, recently, Data-Driven Models (DDMs) are attracting the attention of the industry and academia for their ability to accurately surrogate complex experimental (e.g., EFD) [26], [27], [28], [29] or numerical (e.g., CFD) [1], [14], [16], [17], [19], [21], [22] procedures based on a historical collection of their inputs and outputs, with a function that is computationally expensive to construct but computationally inexpensive to use. Consequently, DDMs can be included directly both in a human-driven optimization loop reducing the computational requirements (i.e., time) between design iterations or developing a fully automated optimization loop requiring minimal human intervention, enabling the exploration of a wider design space [2], [30].

Current approaches to vessel hull resistance optimization rely on a mix between human experience and DDMs [1], [14], [16], [17], [19], [21], [22], [31], [32]. As the first step, human experts define a specific parametrization, i.e., a rich yet synthetic quantitative descriptors of a set of candidate geometries, and parameter ranges, i.e., the geometry design space [1], [14], [16], [17], [19], [21], [22], [31], [32]. For this purpose, several approaches exist in the literature: from Free-Form Deformation (FFD) [1], [16], [17], [19], [21], [22] to B-Splines [14], [22], and model design parameters [31], [32] each one having its strengths and weaknesses (Section II). Once the parametrization and parameter ranges have been defined, a dataset composed of parameters' values (using the selected parametrization) and associated resistance (measured with EFD or estimated using CFD) is built [1], [14], [16], [17], [19], [21], [22], [32]. This process is time, computational, and financially demanding [33], [34], [35]. For this reason, it is necessary to carefully select a minimal number of parameters configuration in the parameters ranges, i.e., a small number of candidate geometries, selected with more or less complex strategies [36], [37] and then perform

the EFD or run the CFD simulations. EFD are seldom used because of the very specific parametrization, and parameter ranges. Moreover, EFD data are seldom shared and available to researchers and practitioners [38] in many cases due to confidentiality issues. In some cases, an already available set of EFD or CFD is available, and it is possible to enrich it with very few new candidate geometries performing EFD or running CFD simulations [39], but to the best of the authors' knowledge, no one in the literature is proposing this approach. Most, if not all, of the work relies just on CFD simulations [14], [16], [17], [19], [21], [22], [32]. Based on the dataset of candidate geometries and their resistance, a DDM-based surrogate of the relationship between the parametrization and the resistance is built, which allows estimating the resistance for a new parameter configuration at a fraction of the time, computational, and financial requirements of the EFD or CFD or both [1], [14], [16], [17], [19], [21], [22], [31], and [32]. The resulting surrogate is then exploited, with different levels of human supervision, by an optimizer to search for the optimal parameters configuration in the parameter range, retrieving then the associated optimized geometry [1], [14], [16], [17], [19], [21], [22], [32]. In practical cases, resistance is one of the different design optimality conditions (e.g., resistance at high and low speed), therefore, multiple optimal solutions are retrieved according to the Pareto front [40]. Figure 1 summarizes the current approach we just described.

The current approach has its limitations. The first one is the need for more or less partial human supervision in geometry parametrization and optimization [1], [14], [16], [17], [19], [21], [22], [31], [32]. In fact, the parametrization needs to satisfy multiple functional requirements: it must be informative enough to allow for the prediction of the resistance and to be homomorphic (i.e., one geometry corresponds to a particular value of the parameters and vice-versa), but it should be synthetic and intelligible enough to allow for interpretation and test (e.g., for physical plausibility of the results) [1], [14], [16], [17], [19], [21], [22], [31], [32]. Moreover, human intervention should also be limited during the optimization phase: the parametrization and the surrogate should be accurate and physically plausible enough to not induce the optimizer into unfeasible, physically implausible, or degenerate solutions [41]. The second limitation is the need for extensive computational efforts (for CFD), costs (for EFD), and time (for both CFD and EFD) needed to build the dataset required in the surrogation [1], [14], [16], [17], [19], [21], [22], [31], [32]. The ideal situation would be to just rely on previous CFD and EFD and not requiring new CFD and EFD for a new design. The last limitation is the limited ability of the approach to work beyond the specific setting (e.g., changes in the family of geometry or extrapolation outside the parameter ranges) as observed in many works [1], [14], [16], [17], [19], [21], [22], [31], [32].

To overcome the limitations discussed above, we propose a novel approach to vessel hull optimization, summarized in Figure 2.

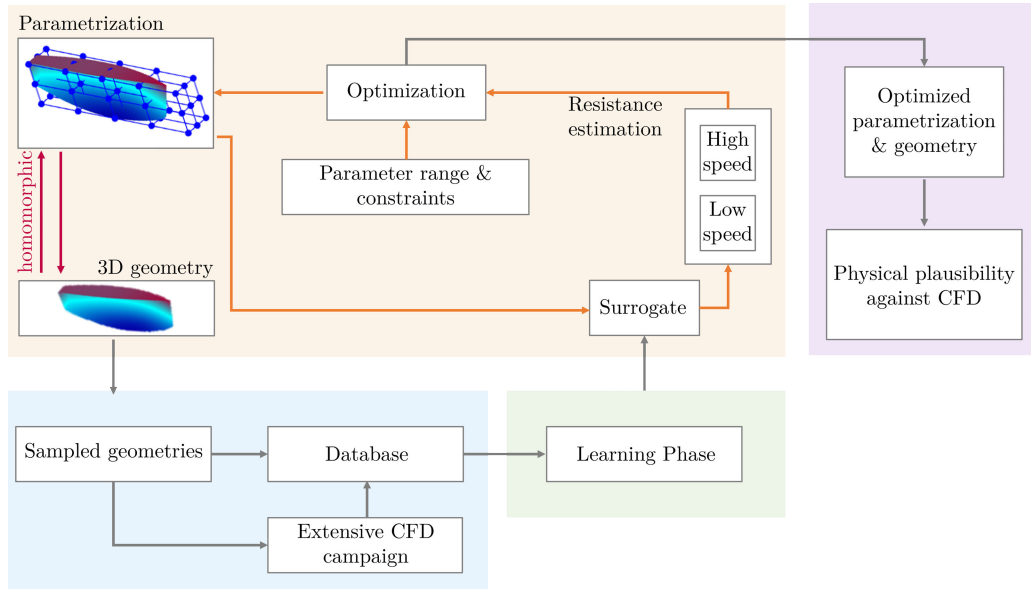


FIGURE 1. Current approach to vessel hull optimization.

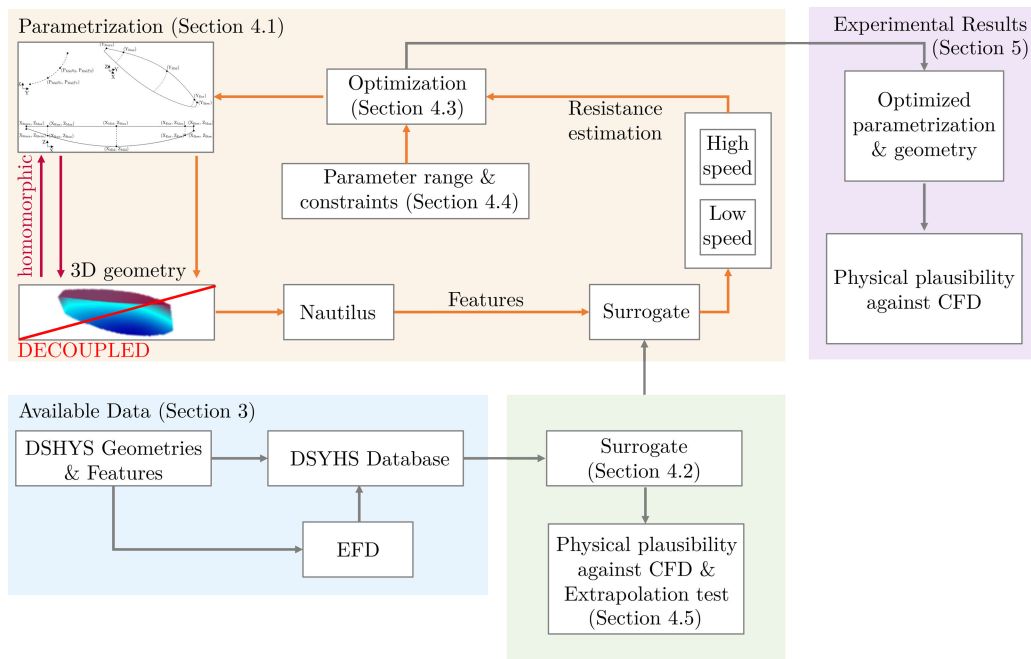


FIGURE 2. Proposed approach to vessel hull optimization.

In previous work [42], we explored the viability of the novel approach in a constrained context, showcasing its effectiveness in a singular scenario. To build on these findings, this study enhances the initial methodology by providing a comprehensive suite of metrics for both development and testing in varied scenarios. In particular, we have broadened the range of algorithms used to develop surrogates, refined both quantitative and qualitative evaluation metrics for surrogates and geometries, and conducted validations across diverse scenarios. This approach offers a more thorough insight into the performance capabilities of the

proposed method and presents a holistic description of the implementation for enhanced repeatability.

As the first step, we propose a parametrization approach able to cover a large set of geometries (i.e., parent hulls) and not just a variation of a particular parent hull. More specifically, our parametrization is a homomorphy not only able to well represent the entire Delft Systematic Yacht Hull Series (DSYHS) (composed of 6 parent hulls) but, as described later, is also able to perform well beyond the DSYHS (i.e., extrapolate). While requiring some human intervention, this step has the capability to minimize it.

In fact, this parametrization can be used for all designs around the 6 parent hulls of the DSYHS, namely, the parametrization step should not be performed every time we change the parent hull as it happens now. This approach paves the way toward more general homomorphic parametrization able to cope with the largest possible sets of parent hulls, allowing us to easily plug them into our pipeline, further decreasing the need for human intervention.

Then, in order to further minimize the human intervention in the parametrization phase, we decoupled the parametrization exploited to define the geometry and to define the optimization parameters from the features necessary to predict the hull resistance based on DDMs. In particular, from a hull geometry defined by a particular configuration of parameters, we exploit the Nautilus code¹ which is able to automatically extract a series of features able to cover and extrapolate over multiple parent hulls while being informative enough to allow for effective and physically plausible predictions of the resistance associated with a particular hull [26], [27], [28], [29]. This decoupling is a fundamental and key contribution to our approach. In fact, the features extracted by Nautilus should not meet the requirement of the geometry parametrization to be homomorphic. This, on one hand, facilitates the ability to create a rich and informative features set that can be used to predict the hull resistance via DDMs without any human intervention (using Nautilus). On the other hand, the homomorphic parametrization of the geometry just needs to focus on the parameters to optimize during the optimization step reducing its complexity and minimizing the human intervention in those cases when new parent hulls need to be covered. This decoupling reduces the original complex and constrained problem into two simpler ones.

Thanks to the decoupled approach to the parametrization, which is able to cover multiple parent hulls [42], we can train DDMs based on the already available EFD of the DSYHS requiring no additional EFD or CFD. However, CFD has been used to check the physical plausibility of the trained DDMs in both synthetic extrapolating scenarios inside the DSYHS and also with a more realistic test outside the DSYHS. For the first case, we defined three, increasingly challenging, extrapolation cases by removing part of the EFD during the DDMs training phase and using those data for testing purposes

- Leave One Velocity Out (LOVO) where we remove all the EFD corresponding to a particular velocity;
- Leave One Geometry Out (LOGO) where we remove all the EFD corresponding to a particular geometry (variation of a particular hull);
- Leave One Series Out (LOSO) where we remove all the EFD corresponding to a particular series (all variations of a particular parent hull).

For the more realistic test outside the DSYHS, we rely on all but one series belonging to DSYHS to train the DDM, and

then we tested it with variations of a particular parent hull that was not used to train the DDM and explore geometric parameters $\delta\%$ larger than the ones covered by the DSYHS.

The proposed surrogate (tested in terms of different extrapolating scenarios and physical plausibility against CFD) is exploited (with minimal levels of human supervision to define the parameters range and constraints) by an optimizer (chosen according to the best options in the literature) to search for the optimal parameters configuration, and retrieve the associated optimized geometry. In particular, we will search for the Pareto front in terms of resistance at high and low speeds. Furthermore, we show that it is possible to optimize the hull resistance by exploring geometric parameters beyond the boundaries of the DSYHS and validating the results via state-of-the-art CFD.

The rest of the paper is organized as follows. Section II reviews the relevant related works; Section III describes the available data; Section IV outlines the proposed methodology; Section V contains the results; and finally, Section VI concludes the work.

II. RELATED WORKS

In this section, we review the relevant related works describing the adopted parametrization, data sources, data validation, surrogates, optimization strategies, obtained results, and the respective validations for each work.

In [14] the authors considered the design optimization of one of the Series 60 hulls. Authors leveraged 7 design variables to parametrize the hull, and an initial population of 210 geometries are evaluated using a Steady Ship Flow solver based on the Neumann-Michell theory. The Steady Ship Flow model was validated against data coming from the literature on the topic. A Radial Basis Function based surrogate of the resistance was constructed and leveraged within Multi-Objective Artificial Bee Colony based optimization framework. Statistical validation for the surrogate was not reported, but a visual representation of the results, i.e., scatter plots of the real versus predicted resistance showed a good agreement between the surrogate and the original model. The optimization was constrained by a 1% change in the displacement. The results of the optimization found an 8% reduction in total resistance at a Froude number² of 0.27 and were validated using a Reynolds Averaged Navier Stokes based CFD tool (which is a high fidelity approach).

In [1] the authors relied on CFD and EFD to optimize an offshore aquaculture vessel. The ship was parametrized using FFD with 9 design variables, and 300 geometries were sampled and evaluated using a CFD model validated against EFD data. The CFD had a good agreement with the EFD and showed a maximum error of 6.7%. The authors proposed a Support Vector Regression based surrogate coupled with the NSGA-II optimizer to minimize total resistance and wake

²The Froude number is commonly used in naval architecture to represent the ratio between the inertial and gravitational forces, and is proportional to the velocity [43].

¹<https://github.com/mai-lab-tud/nautilus>

flow at the design speed. The results showed reductions of 1.6% and 18% for total resistance and wake flow, respectively. A model of the optimal hull was built and tested via EFD to validate the reduction in total resistance and wake flow.

In [16] the authors considered the design optimization of a bulbous bow. The bow was parametrized based on FFD with 6 design parameters to control the protrusion and immersion of the bulb. A small number of geometries (25) were evaluated using high fidelity CFD considering 5 Froude numbers (i.e., 0.294, 0.312, 0.331, 0.349, 0.367, and 0.386). The CFD was validated with two levels of mesh coarsening (0.8M and 2.8M cells) and compared to EFD which showed a maximum deviation of 5.6%. A Krigging based surrogate was developed and coupled with a NSGA-II optimizer. The 5 best geometries (each corresponding to different local minima) were validated against the CFD and the results showed that the performance of 3 geometries were unsatisfactory while the remaining 2 showing a reduction in the total resistance of approximately $6 \div 7\%$.

In [32] the authors parametrized a twin-skeg fishing vessel based on 6 dimensionless design variables. The CFD approach was numerically validated through a mesh-coarsening procedure with 3 levels of grid refinement and 54 simulations were performed to construct a dataset. A Krigging based surrogate showed a Coefficient of Determination of ~ 0.95 . Coupling the latter with a NSGA-II optimizer, the authors minimized the total resistance at 4 different velocities between 9 and 12 knots. The optimal model corresponded to a reduction of 5.6% in the total resistance at the design speed of 11.3 knots which was later validated using the CFD model.

In [17] the authors addressed the optimization of the KCS vessel parametrized based on FFD and 6 design variables. A high fidelity CFD model was constructed to evaluate the performance of 120 geometries at the design speed (Froude number was fixed to 0.26). A mesh coarsening procedure was exploited to validate the mesh (1.5M cells) along with a comparison to EFD coming from the literature. A Response Surface Model based surrogate was constructed and statistically validated with Leave One Out showing a Coefficient of Determination of 0.97 and Root Mean Square Deviation of 0.05N. The NSGA-II based optimization framework was employed to minimize the total resistance subject to a 1% change in displacement. The results, validated with CFD model, showed a reduction of 0.32N in the total drag.

In [19] the authors considered the optimization of the KCS vessel at 2 speeds. The vessel was parametrized based on FFD using 5 design variables. A Neumann-Michell coupled with Reynold Averaged Navier Stokes CFD based approach, validated on EFD, was exploited to generate an initial population of 40 geometries evaluated at 2 speeds. A Gaussian Progress Regression based surrogate was developed, but no statistical validation was reported, nevertheless scatter plots of real versus predicted resistance

show a good agreement between the surrogate and the Neumann-Michell coupled with CFD based approach. The authors leverage a NSGA-II optimizer to minimize the resistance, and the results showed a 9.24% reduction at a Froude number of 0.26 (corresponding to the vessel design speed) and a 4.99% reduction at a Froude number of 0.2.

In [21] the authors parametrized a Wigleyship based on FFD using 2 design variables. A panel-based CFD approach was exploited to evaluate the wave making coefficient at a Froude number of 0.35 and the model was validated against CFD data coming from the literature. A Deep Belief Neural Network based surrogate showed a Coefficient of Determination of ~ 1 . Coupling the latter with a Quadratic Lagrangian based Non-Linear Programming optimizer, the authors minimized the wave making coefficient at Froude numbers between 0.28 and 0.36. The optimal model corresponded to a reduction of 12% in the wave making coefficient which was later validated using a CFD model showing an uncertainty of 5%.

In [31] the authors leveraged a parametric model consisting of 6 design variables to optimize a catamaran. The authors evaluated 2000 geometries using a panel-based CFD model validated against high fidelity CFD results coming from the literature. This study included a notably higher number of geometries with respect to other referenced works due to their use of low fidelity CFD which is of course less computationally expensive. Gaussian Process, Support Vector, and Multi Adaptive Splines based surrogates were developed and statistically validated showing a Coefficients of Prognosis of 0.57, 0.73, and 0.81 respectively. For this reason, the authors exploited the Multi Adaptive Splines based surrogate coupled with the NSGA-II optimizer to minimize the vessel resistance at 21, 23, 25, 27 and 30 knots. Two optimal hull designs were validated against the low fidelity CFD and a high fidelity CFD with 2.1M cells. For one of the two hulls, the panel-based CFD predicted a decrease in resistance of $1.32 \div 1.44\%$ but the high fidelity CFD actually showed an increase in resistance of 1.2%. For the other hull the panel-based CFD predicted a reduction of $0.91 \div 1.57\%$ which was in agreement with the high fidelity CFD showing a reduction of 1.1%.

In [22] the authors leveraged a FFD and spline based parametrization with 5 design variables for the DTMB-5415 hull. The authors combined the Neumann-Michell and the Reynolds Averaged Navier Stokes CFD models to predict the coefficient of resistance by evaluating 50 geometries using Neumann-Michell and 30 geometries using the Reynolds Averaged Navier Stokes CFD models. The mesh was validated through a coarsening procedure. The Neumann-Michell and the Reynolds Averaged Navier Stokes CFD were validated on experimental data coming from the literature. A Krigging based surrogate was proposed and statistically validated showing an Average Absolute Error of 0.29, a Maximum Absolute Error of 1.93, and a Root Mean Square Error of 0.45. The authors demonstrated that the surrogate accuracy increased with the addition of

TABLE 1. Summary of the reviewed related works reported in Section II describing the adopted parametrization, data sources, data validation, surrogate and optimization strategies, obtained results, and the respective validations for each one of them.

Ref.	Parametrization	Data Sources	Data Validation	Surrogate	Surrogate Validation	Optimization	Results	Optimization Validation
[14]	7 design variables	210 samples generated with a Steady Ship Flow solver based on the Neumann-Michell theory	Against data coming from the literature on the topic	Radial Basis Function	-	Multi-Objective Artificial Bee Colony to minimize total resistance subject to the displacement constrained by $\pm 1\%$	8% reduction in total resistance	The optimal design was validated with a Reynolds Averaged Navier Stokes based CFD
[11]	FFD with 9 design variables	300 samples generated with high fidelity CFD	Against EFD data, maximum error of 6.7%	Support Vector Regression	-	A NSGA-II based framework to minimize the total resistance and the wake coefficient, subject to constraints applied to the parameters	Reductions of 1.6% and 18% for total resistance and wake flow respectively	The optimal design was validated against EFD data
[16]	FFD with 6 design variables	25 geometries were evaluated using high fidelity CFD considering 5 Froude numbers (i.e., 0.294, 0.312, 0.331, 0.349, 0.367, and 0.386)	Two levels of mesh coarsening (0.8M and 2.8M cells) and compared to EFD which showed a maximum deviation of 5.6%	Krigging	-	A NSGA-II based framework to minimize the total resistance, subject to constraints applied to the parameters	Reductions of 5÷6% in the total resistance	The 5 best geometries were validated against the CFD: 3 showed unsatisfactory and 2 showed a reduction in the total resistance of approximately 6÷7%
[32]	6 design variables	54 geometries were evaluated using high fidelity CFD considering 4 velocities	Three levels of mesh coarsening (0.07M, 0.09M and 0.11M cells)	Krigging	-	A NSGA-II based framework to minimize the total resistance, subject to constraints applied to the parameters and submerged volume	Reductions of 5.6% in the total resistance	The best geometry was validated against the CFD model
[17]	FFD with 6 design variables	120 samples from high fidelity CFD	Mesh coarsening to validate the mesh (1.5M cells) along with a comparison to EFD coming from the literature	Response Surface Model	Leave One Out (Coefficient of Determination equal to 0.97) and Root Mean Square Deviation of 0.05N	A Genetic Algorithm based framework to minimize total resistance, subject to constraints applied to the parameters	A reduction of 0.32N in total resistance	The optimal design was validated with high fidelity CFD
[19]	FFD with 5 design variables	40 geometries evaluated at 2 speeds for a total of 80 samples generated with a Neumann-Michell coupled with CFD based approach	Against EFD data	Gaussian Process Regression	Good agreement according to several scatter plot	A NSGA-II based framework to minimize total resistance at 2 speeds, subject to the displacement and surface area constrained by $\pm 1\%$	A reduction of 9.24% in total resistance at a Froude number of 0.26 & 4.99% at a Froude number of 0.2	-
[21]	FFD with 2 design variables	100 geometries evaluated using low fidelity CFD	CFD data coming from the literature	Deep Belief Network developed	Coefficient of Determination of ~ 1	Quadratic Lagrangian based Non-Linear Programming framework to optimize the wave making coefficient, subject to constraints applied to the parameters	A reductions of 12%	The optimal hull was validated with high fidelity CFD showing an uncertainty of 5%
[31]	6 design variables	2000 samples from low fidelity CFD	High fidelity CFD results coming from the literature	Gaussian Process, Support Vector, Multi Adaptive Spline Regressors	Coefficients of Prognosis of 0.57, 0.73, and 0.81 respectively	A NSGA-II based framework to minimize total resistance at 5 speeds	1.2% increase in resistance for one hull & 1.1% reduction for the second hull	The optimal hulls were validated with a high fidelity CFD
[22]	FFD coupled with spline approach with 5 design variables	50 geometries evaluated with a Neumann-Michell model & 30 geometries evaluated with a Reynolds Averaged Navier Stokes model	The mesh was validated through a coarsening procedure. The Neumann-Michell and the Reynolds Averaged Navier Stokes CFD were validated on experimental data coming from the literature	Krigging (Mean Absolute Error of 1.93)	Average Absolute Error of 0.29, a Maximum Absolute Error of 1.93, and a Root Mean Square Error of 0.45.	Genetic Algorithm based framework to minimize total resistance a Froude number of 0.28	A reduction of 2÷5% in total resistance	The optimal hull was validated with high fidelity CFD

high fidelity samples. The optimization was carried out using a Genetic Algorithm to minimize the resistance at a Froude number of 0.28 showing a reduction in resistance of 2÷5%.

For the sake of completeness, the review's summary is also reported in Table 1.

From the literature review, it is possible to identify several areas of interest for future research to address the current gaps in methodology.

Specifically, the current approach utilizes FFD and parametric model based parametrizations, with a focus on building a design space around a specific parent hull. Additionally, data generation and validation is completed by the use of both CFD and EFD, which suffers from significant computational and time demands. There is a diverse use of surrogates, however, areas of surrogate validation are not well defined, often lacking statistical validation, relying instead on qualitative agreements from scatter plots. This is crucial as it was observed that the optimizer could be induced into false minima that were not physically plausible during development.

In contrast, the proposed approach addresses all of these challenges simultaneously. In particular, the shape decoupled methodology allows for the use of historical data in the surrogate training phase and the robust statistical validation ensures the surrogates are physically plausible. The comparison between current approaches in the literature and the proposed approach are summarized in Table 2.

III. AVAILABLE DATA

In this section we will describe the data that we will exploit in this study. In particular, we leverage the DSYHS database [44] (available upon request to the Delft University of Technology Ship Hydromechanics Laboratory³) which has been used in a number of works [26], [27], [28], [29], [45].

In [11] the authors present the original series of the DSYHS which included 22 systematically varied sailing yacht hulls, alongside a polynomial expression they developed to determine the residual hull resistance in terms of the hull geometry, over a range of Froude numbers. In the successive years many more experiments were added to the DSYHS database

³<https://dsyhs.tudelft.nl>

TABLE 2. Comparison between the current approaches (see Table 1) and the proposed one describing the adopted parametrization, data sources, data validation, surrogate and optimization strategies, obtained results, and the respective validations for each one of them.

App.	Parametrization	Data Sources	Data Validation	Surrogate	Surrogate Validation	Optimization	Results	Optimization Validation
Current	EFD or Parametric models with 2 ÷ 9 design variables based around a particular parent hull.	Use of CFD and EFD for geometry generation and validation. Sample sizes from 25 ÷ 2000 geometries depending on fidelity.	CFD models validated against EFD data or high-fidelity CFD, with reported errors up to 6.7%.	Radial Basis Function, Support Vector Regression, Kriging, Gaussian Process Regression, Deep Belief Networks.	Not always reported. Some mention good agreement based on scatter plots but not focused on physical plausibility.	Multi-Objective Artificial Bee Colony, NSGA-II, Genetic Algorithms to minimize resistance, or optimize wake flow.	Reductions in total resistance ranging from 0.91% ÷ 12%.	Validation performed using high-fidelity CFD models or EFD but showed a lack of physical plausibility in some instances.
Proposed	Section IV-C Homomorphic parametrization for the entire DSYHS and beyond, designed to minimize human intervention with reparametrizing.	Section III Existing EFD data, eliminating the need for new EFD or CFD.	— EFD data only.	Section IV-A DDMs trained on existing EFD data with automatic feature extraction via Nautilus.	Section IV-B Validated for physical plausibility in different extrapolation scenarios against state-of-the-art CFD models.	Section IV-D Optimizer chosen based on best literature practices, aiming to find the Pareto front for resistance.	Section V Results for extrapolation scenarios when exploring geometric parameters beyond the DSYHS.	Section IV-E High-fidelity CFD used for validation.

TABLE 3. Geometric boundaries of each series of the DSYHS.

Series	Geometries	Length [m]	Breadth [m]	Depth [m]	Volume $\times 10^{-3} [m^3]$	Draft [m]
S_1	22	1.94 ÷ 2.30	0.51 ÷ 0.68	0.29 ÷ 0.36	37.60 ÷ 37.61	0.13 ÷ 0.15
S_2	6	2.31 ÷ 2.35	0.50 ÷ 0.66	0.23 ÷ 0.35	37.53 ÷ 37.60	0.11 ÷ 0.15
S_3	11	2.28 ÷ 2.40	0.48 ÷ 0.80	0.20 ÷ 0.28	37.53 ÷ 37.53	0.11 ÷ 0.11
S_4	9	2.09 ÷ 2.38	0.37 ÷ 0.74	0.27 ÷ 0.34	30.92 ÷ 37.53	0.11 ÷ 0.12
S_6	3	2.42 ÷ 2.43	0.66 ÷ 0.66	0.30 ÷ 0.35	30.92 ÷ 30.92	0.12 ÷ 0.12
S_7	3	2.51 ÷ 2.57	0.38 ÷ 0.45	0.21 ÷ 0.33	30.92 ÷ 30.92	0.12 ÷ 0.12

TABLE 4. Hydrostatic coefficients provided in the DSYHS to describe the geometries.

Parameter	Symbol
Length of waterline	L_{wl}
Breadth of waterline	B_{wl}
Draft of canoe body	T
Volume of canoe body	∇
Longitudinal center of buoyancy	LCB
Longitudinal center of flotation	LCF
Area of waterplane	A_w
Area of cross-section	A_x
Wetted surface area of canoe body	S
Block coefficient	C_b
Midship area coefficient	C_m
Prismatic coefficient	C_p
Waterplane area coefficient	C_w

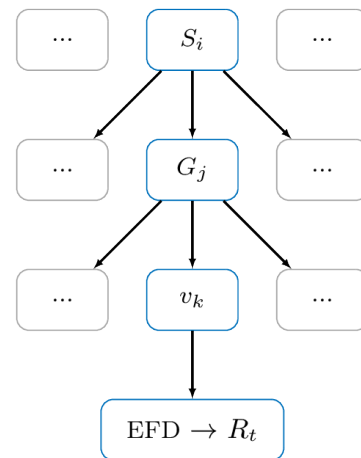
and now, to the best of the authors' knowledge, the DSYHS database is currently the largest collection of sailing yacht EFD in the world.

The current DSYHS database contains the hull collections for Series 1 ÷ 7 ($S_1 \div S_7$) where S_5 does not exist in the database. The 6 series, composed of parent hulls and their derivatives, are in model scale (which is the scale which the experiments were performed at) and span approximate lengths of 2.100 ÷ 2.500m, widths of 0.440 ÷ 0.660m, and depths of 0.270 ÷ 0.350m. Table 3 shows the geometric boundaries of each series of the DSYHS.

From these 6 Series, namely the parent hulls, 54 different geometries G have been derived. For each geometry, the total resistance R_t over a range of speeds v had been retrieved via EFD. A visual representation of this description is reported in Figure 3. The total number of EFD in the DSYHS dataset is 702.

The 54 geometries contained in the DSYHS are described through the use of hydrostatic coefficients common for naval architecture applications, see Table 4 for details.

Note that, for a general geometry, the parameters reported in Table 4 can be easily retrieved with Nautilus¹.

**FIGURE 3.** The database contains the hull collections for Series 1 ÷ 7 ($S_1 \div S_7$) where S_5 does not exist in the database. From these 6 Series, namely the parent hulls, 54 different geometries G have been derived. For each geometry, the total resistance R_t over a range of speeds v had been retrieved via EFD.

IV. METHODOLOGY

In this section, we will deepen the description of the methodology we propose starting from the schema presented in Figure 2.

In particular, Section III already focused on the available data, while the following aspects of the methodology will be the subjects of this section:

- the development of the surrogate to estimate R_t based on the parameters reported in Table 4 that can also be retrieved with Nautilus¹ for any hull geometry (Section IV-A);
- the validation in different extrapolating scenarios and the physical plausibility against CFD of the surrogate (Section IV-B);
- the homomorphic parametrization of the hull and the parameters range and constraints generating the parameter space (Section IV-C);
- the optimization framework which searches in the homomorphic parameters space simultaneously optimizing R_t for both a high v^{High} and low v^{Low} estimated with the surrogate (Section IV-D);

- the verification of the physical plausibility against CFD of the geometries on the Pareto front generated by the optimizer (Section IV-E).

Note that, with the aid of the proposed decoupling strategy between the parametrization exploited by the optimizer and the one exploited by the surrogate, given a point in the homomorphic parameters space it is possible to extract the input of the surrogate with Nautilus¹ and estimate the R_t for v^{High} and v^{Low} with minimal computational requirements making the optimization fast and cheap to perform.

A. SURROGATE DEVELOPMENT

The problem of predicting R_t based on the parameters reported in Table 4 and the velocity v can be mapped to a typical regression problem by Machine Learning [46], [47].

The No-Free-Lunch Theorem [48] ensures us that, in order to find the best algorithm for a particular application, it is necessary to test multiple algorithms. In our case, we will test 4 state-of-the-art algorithms⁴ [49], [50]: Random Forests (RF) [51], [52], XGBoost [53], Kernel Ridge Regression (KRR) [47], and the Extreme Learning Machine (ELM) [54], [55] namely a Single Layered Neural Network [56], [57] where the weights of the first layers have been randomly set reducing the computational burden of the training phase with minimal, if not absent, effect on accuracy.

In RF we need to tune the number of features to randomly sample from the whole features during each node of each tree creation n_f and the maximum number of elements in each leaf of each tree n_l . As RF performance improves by increasing the number of trees n_t , we set it to 1000 as a reasonably large number yet computationally tractable.

In XGBoost, we need to tune the learning rate of the gradient l_r , the max depth of each tree n_d , the minimum loss reduction m_l , the number of points to randomly sample from the whole training set for each tree creation n_b , and the number of features to randomly sample from the whole training set during the creation of each node for each tree n_f .

In KRR we chose to rely on the Gaussian kernel for the reason described in [58], and then the regularisation hyperparameter λ and the kernel coefficient γ need to be tuned.

In ELM, we use the sigmoid activation function in the hidden layer and the linear activation in the output layer. Then we need to tune the number of hidden neurons h_l and then the regularisation hyperparameter λ on the weights of the last layer.

The summary of these hyperparameters with the associated search space is reported in Table 5.

Note that, the selection of the best performing algorithm and the best hyperparameters, will depend on the scenario under consideration and on two different metrics, namely accuracy and computational requirements (see Section IV-B).

TABLE 5. Hyperparameters and hyperparameters search space for all algorithms tested in this work, $d = 13$ denotes the number of features in the dataset (see Table 4).

Algorithm	Hyperparameters
RF	$n_f : \{d^{0.33}, d^{0.5}, d^{0.75}\}$ $n_l : \{1, 3, 5, 10\}$ $n_t : \{1000\}$
XGBoost	$l_r : \{0.01, 0.02, 0.03, 0.04, 0.05\}$ $n_d : \{3, 5, 10\}$ $m_l : \{0, 0.1, 0.2\}$ $n_b : \{0.6n, 0.8n, 1n\}$ $n_f : \{0.5d, 0.8d, 1d\}$
KRR	$\lambda : \{10^{-6}, 10^{-5.8}, \dots, 10^3\}$ $\gamma : \{10^{-6}, 10^{-5.8}, \dots, 10^3\}$
ELM	$h_l : \{2^5, 2^6, \dots, 2^{16}\}$ $\lambda : \{10^{-6}, 10^{-5.8}, \dots, 10^3\}$

The performance, in terms of accuracy, will be measured in accordance with different metrics: three quantitative (the Mean Absolute Error - MAE, the Mean Absolute Percentage Error - MAPE, and the Pearson Product-Moment Correlation Coefficient - PPMCC) [59] and one qualitative (the scatter plot actual versus predicted value) [60].

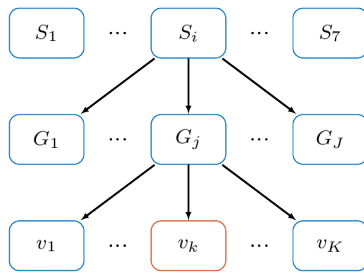
The performance, in terms of computational requirements, will be measured by means of time to build the model (Training Time) and time to make a prediction (Test Time). Since our surrogate will be leveraged in the optimization phase (see Section IV-D), the most important computational metric is the Test Time.

B. SURROGATE VALIDATION AND PHYSICAL PLAUSIBILITY

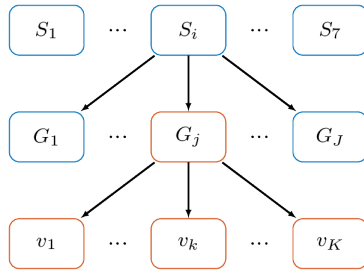
In our work, we will study three different extrapolating scenarios based on the intrinsic hierarchy of the dataset. This will allow us to understand the extrapolation ability and the robustness of the different models described in the previous section (see Figure 4 for a visual representation):

- LOVO: where we remove all the EFD corresponding to a particular velocity. Since the EFD, for each geometry and each series, has been performed at different speeds, we create an histogram of the velocities with 16 bins. For the sake of replicability, one can find the final binning (with lower v_l and upper v_u bounds) reported in Table 6. The LOVO scenario, then, is actually leaving out all the EFD following in one of these bins. The scope of this scenario is to test the extrapolation ability of the model in terms of velocity, namely to estimate the resistance at a velocity never observed before in the dataset;
- LOGO: where we remove all the EFD corresponding to a particular geometry (variation of a particular hull). The scope of this scenario is to test the extrapolation ability of the model in terms of geometry, namely to estimate the resistance of a geometry never observed before in the dataset;
- LOSO: where we remove all the EFD corresponding to a particular series (all variations of a particular parent hull). The scope of this scenario is to test the extrapolation ability of the model in terms of series,

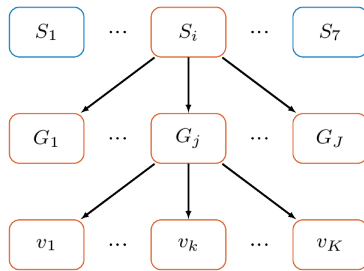
⁴Results in Kaggle www.kaggle.com, the most popular Machine Learning competition website, shows that these algorithms are the top winners.



(a) LOVO.



(b) LOGO.



(c) LOSO.

FIGURE 4. Visual representation of the three different extrapolating scenarios we investigated in this work based on the intrinsic hierarchy of the dataset. In particular we highlighted data hidden from the training phase and exploited just for testing purposes in orange.

TABLE 6. Histogram of the velocities with 16 bins for the DSYHS EFD.

Bin	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
v_l	0.00	0.30	0.62	0.94	1.26	1.58	1.90	2.22	2.54	2.86	3.18	3.50	3.82	4.14	4.46	4.78
v_u	0.30	0.62	0.94	1.26	1.58	1.90	2.22	2.54	2.86	3.18	3.50	3.82	4.14	4.46	4.78	5.10

namely to estimate the resistance for a series never observed before in the dataset.

Note that the LOSO scenario is, in our work, the most interesting and useful one in practical applications. In fact, in practice, what we want to do is to generate geometry for a new, previously unexplored series, and this is precisely the scope of the LOSO scenario: we assume to have developed a few series, and we try to infer something for a news series that was previously unexplored.

What remains to be addressed is how to tune the hyperparameters of the different Machine Learning algorithms that we tested to generate the surrogate (see Section IV-A) and how to assess their final performance [61].

For what concerns the last point, the answer is easy. Based on the different scenarios (LOVO, LOGO, and LOSO) we have to split the data in Training \mathcal{D}_n and Test \mathcal{T}_t sets using the principle of the different extrapolating scenarios. For example, in the LOVO scenario, we put all the EFD corresponding to one of the histogram bins in \mathcal{T}_t while the remaining ones are kept in the \mathcal{D}_n . Then we can use \mathcal{D}_n to both train the model and select the associated best hyperparameters and use \mathcal{T}_t to assess the performance of the final model. Repeating multiple times, this procedure will give us the average performance in the different scenarios.

Instead, for tuning the hyperparameters of the different Machine Learning algorithms, we proceeded as follows. We took \mathcal{D}_n and split it into Learning \mathcal{L}_l and Validation \mathcal{V}_v sets considering the very same extrapolating scenario that we use for assessing the final performance. Then we train each model with \mathcal{L}_l with many different hyperparameters configurations and measure its performance on \mathcal{V}_v according to the MAE. Then we repeated the experiment multiple times and selected the hyperparameters' configuration which gives the best average MAE on the validation sets. Finally, we retrained the model with the selected best configuration of the hyperparameters on the whole \mathcal{D}_n which is the model that will be used for testing purposes (see the previous paragraph).

To ensure the physical plausibility of the proposed surrogate, we leveraged a state-of-the-art CFD model. For the DSYHS, EFD have been carried out by means of a large experimental campaign carried out at the Delft University of Technology towing tank⁵ and for this reason, they possess some level of uncertainty that cannot be removed. Therefore, to measure the quality of our surrogate we need to compare its performance against a baseline which, in our case, is a state-of-the-art CFD model.

Unfortunately, the CFD model is too computationally expensive to run for all the geometries and velocities in the databases. For this reason, we will compare our surrogate on a subset of them. In particular, we will consider the most challenging scenario, i.e., the LOSO, and we will perform the comparison between the CFD and the proposed surrogate models on the series which exhibit the largest deviation between the surrogate and the EFD results.

For the CFD model, the mesh generation, the computation of the solution, and the post-processing of results was carried out in Star CCM+⁶ which is a state-of-the-art commercial CFD package. The simulation domain was created to satisfy the following constraints: the depth under the vessel was greater than twice the draft, the length of the domain after the vessel was longer than twice the length of the vessel, and the width of the domain was 50% larger than the length of the vessel. To reduce the computational demand of the simulation, the hull was divided symmetrically along the longitudinal axis and only half of the problem was simulated

⁵www.tudelft.nl/3me/over/afdelingen/maritime-and-transport-technology/research/ship-hydraulics/facilities/towing-tank-no-1

⁶www.plm.automation.siemens.com/global/en/products/simcenter/STAR-CCM.html

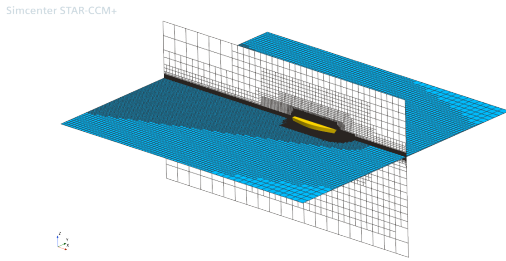


FIGURE 5. The mesh exploited for the CFD simulations with the Star CCM+⁶ package. The mesh included a surface mesh refinement on the vessel hull and on the boundaries of the domain in addition to volume mesh refinements around the hull, wake, and free surface.

to assess the hydrodynamic performance. The CFD model is a finite volume based viscous RANS solver which can compute the hull resistance in various calm-water conditions by solving the underlying partial differential equations. For the problem at hand, a RANS solver was used considering the $k - \omega$ shear-stress turbulence model with wall functions [62]. The boundaries of the domain were set so the symmetric and far-field boundaries were considered as symmetry planes. The top, bottom, and inlet boundaries were considered as velocity inlets while the outlet boundary was considered as a pressure outlet. The volume of fluid technique was used to establish a free surface in the solution and solving the underlying equations with the volume fraction of both water and air [63]. To find the solution of the hull resistance, the vessel was simulated using the dynamic fluid body interaction module in Star CCM+⁶ with two degrees of freedom (sink and trim), which is in line with the experimental campaign outlined in [11], [12], and [13]. The simulation was set-up with a time-step of 0.001s and the behavior of the vessel simulated for a period of 60s. The solution of the simulation was then taken as the time averaged response over this period. The described CFD simulation was validated against the original EFD results for a number of geometries to ensure it could be used for the physical plausibility of the surrogate. A mesh coarsening procedure was carried out with $3 \cdot 10^5$, $9 \cdot 10^5$, and $3 \cdot 10^6$ cells respectively to ensure there was grid independence. Results using the highest fidelity mesh with $3 \cdot 10^6$ cells are presented in Section V. Figure 5 shows the exploited mesh for the CFD simulations with the Star CCM+⁶ package. The mesh included a surface mesh refinement on the vessel hull and on the boundaries of the domain in addition to volume mesh refinements around the hull, wake, and free surface.

C. HULL PARAMETRIZATION AND PARAMETERS RANGE

In this section, we will describe the adopted homomorphic parametrization, together with the associated parameter range, that will be leveraged during the optimization phase (see Section IV-A) to search for the best hull, i.e., the hull that will exhibit the best R_t at v^{Low} and v^{High} . It is worth noting that this parametrization is decoupled from the one exploited in the definition of the surrogate (see Section IV-D) as described in the introduction (see Figure 2).

In particular, a parametric model for a sailing yacht hull [42], [64], [65], [66], [67], [68], [69], [70] was developed with the Siemens NX⁷ software leveraging on 32 parameters. The full list of parameters together with their description is reported in Table 7 and visualized in Figure 6.

The 32 parameters govern the hull geometry through the use of B-Spline curves [71], which in turn, drive the design of the yacht hull surface inside the parametric model. The parametrization is directly related to control points on the B-Spline curves which allows the parameters to be modified independently and ensures the desired homomorphic properties. Geometric constraints were imposed on the model to ensure G0 (positional) and G1 (tangential) continuity at the intersection between adjacent splines to assist in producing feasible designs. Additionally, the G2 (curvature) continuity was also applied to ensure a smooth surface was retrieved from the model [70]. Figure 6 includes: an example cross-section of the mid section (top left), an isometric view of the parametric hull (top right), and a planar view in the xz plane of the parametric hull (bottom). Parameters denoted with an x or z define features in the xz plane and parameters denoted with y define features in the xy plane. The parameters preceded by P refer to the B-spline control points in the yz plane of each section.

For what concerns the parameters ranges, they have been designed following this principle. First, for each parameter, we search for the minimum and maximum value in a specific series S_i , i.e., the series that we want to optimize (see Section IV-D). Then we increased that range by $\delta\%$. This extrapolation is especially useful because, in practice, we want to be able to generate a geometry for a new, previously unexplored series rather than restrict ourselves to preexisting designs. In the experiments, we will show that the $\delta = 30\%$ is the limit threshold beyond which the surrogate starts to induce the optimizer into degenerate solutions. The parameter ranges extracted from the original 54 hulls belonging to the DSHYS database are reported in Table 8.

The proposed homomorphic parametrization does not succumb to the limitations of the current approaches, i.e., the need to re-parametrize each parent geometry, and is able to cover the whole DSHYS database and beyond (i.e., up to $\delta\%$ of the DSHYS).

D. OPTIMIZATION FRAMEWORK

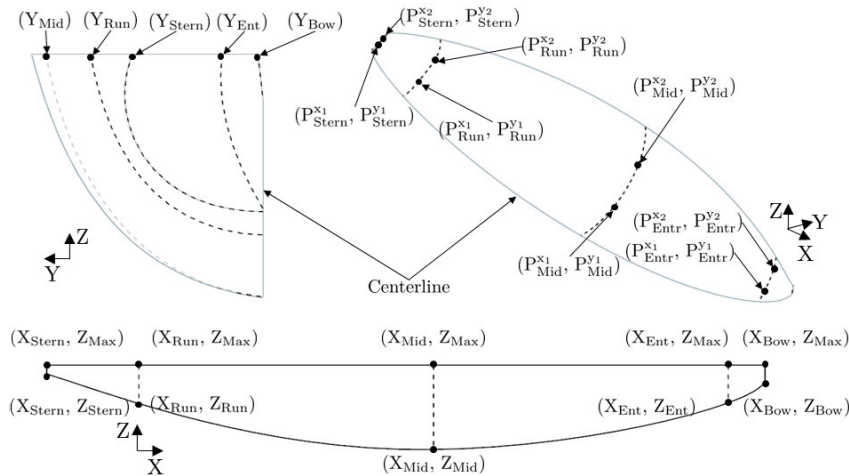
In this section we will present the proposed approach to the search for the best hull in a series S_i , in terms of the best R_t at v^{High} and v^{Low} , in line with the referenced works [16], [31], leveraging the parametrization described in Section IV-C and the surrogate described in Section IV-A.

Since it is not a fair trade-off between the resistance R_t and the volume of the hull ∇ , i.e., having no submerged body will correspond to the case where the resistance is zero, and in line with the original investigations of [11], [12], and [13],

⁷www.plm.automation.siemens.com/global/en/products/nx/

TABLE 7. The 32 parameters together with their description characterizing the adopted homomorphic parametrization for a sailing yacht hull developed with the Siemens NX ⁷.

Symbol	Description	Symbol	Description
X_{Bow}	Length along the center line from the origin to the bow.	P_{Entr}^{y1}	Width from the center line to the first control point of entrance section.
Y_{Bow}	Width from the center line to the bow section at the maximum height.	P_{Entr}^{z1}	Height from Z_{Entr} to the first control point of entrance section.
Z_{Bow}	Height from the origin to the bottom of the bow section.	P_{Entr}^{y2}	Width from the center line to the second control point of entrance section.
Z_{Max}	Height from the origin to maximum height.	P_{Entr}^{z2}	Height from Z_{Entr} to the second control point of the entrance section.
X_{Entr}	Length along the center line from the origin to the entrance section.	P_{Mid}^{y1}	Width from the center line to the first control point of the first mid section.
Y_{Entr}	Width from the center line to the entrance section at the maximum height.	P_{Mid}^{z1}	Height from Z_{Mid} to the first control point of the first mid section.
Z_{Entr}	Height from the origin to the bottom of the entrance section.	P_{Mid}^{y2}	Width from the center line to the first second point of the first mid section.
X_{Mid}	Length along the center line from the origin to the first mid section.	P_{Mid}^{z2}	Height from Z_{Mid} to the second control point of the first mid section.
Y_{Mid}	Width from the center line to the first mid section at the maximum height.	P_{Run}^{y1}	Width from the center line to the first control point of the run section.
Z_{Mid}	Height from the origin to the bottom of the first mid section.	P_{Run}^{z1}	Height from Z_{Run} to the first control point of the run section.
X_{Run}	Length along the center line from the origin to the run section.	P_{Run}^{y2}	Width from the center line to the second control point of the run section.
Y_{Run}	Width from the center line to the run section at the maximum height	P_{Run}^{z2}	Height from Z_{Run} to the second control point of the run section.
Z_{Run}	Height from the origin to the bottom of the run section.	P_{Stern}^{y1}	Width from the center line to the first control point of the stern section.
X_{Stern}	Length along the center line from the origin to the stern section.	P_{Stern}^{z1}	Height from Z_{Stern} to the first control point of the stern section.
Y_{Stern}	Width from the center line to the stern section at the maximum height.	P_{Stern}^{y2}	Width from the center line to the second control point of the stern section.
Z_{Stern}	Height from the origin to the bottom of the stern section.	P_{Stern}^{z2}	Height from Z_{Stern} to the second control point of the stern section.

**FIGURE 6.** A visual representation of the 32 parameters characterizing the adopted homomorphic parametrization for a sailing yacht hull developed with Siemens NX ⁷.

we are concerned about optimizing the relative resistance to the submerged volume (i.e., $\frac{R_t}{\nabla}$).

Again, in line with the referenced works [14], [19], the optimization problem is subject to a constraint to bound the volume according to a lower and upper boundary ∇_l and ∇_u respectively. ∇_l and ∇_u have been set by searching for the minimum and maximum value in a specific series S_i , i.e., the series that we want to optimize, because we aim to optimize the geometry of a hull that fits within a particular series i.e., conforms to the same volume constraints.

At this point, we can formalize our problem as follows

$$\begin{aligned}
 \min_{\mathbf{p}} \quad & \left\{ \frac{R_t(\mathbf{p}, v^{\text{High}})}{\nabla(\mathbf{p})}, \frac{R_t(\mathbf{p}, v^{\text{Low}})}{\nabla(\mathbf{p})} \right\}, \\
 \text{s.t.} \quad & \nabla_l \leq \nabla(\mathbf{p}) \leq \nabla_u, \\
 & \mathbf{p}_l(\delta) \leq \mathbf{p} \leq \mathbf{p}_u(\delta),
 \end{aligned} \tag{1}$$

where \mathbf{p} is the vector of the 32 parameters of the homomorphic parametrization of Table 7, $\mathbf{p}_l(\delta)$ and $\mathbf{p}_u(\delta)$ are their lower and upper bounds of the parameters as a function of δ , $\nabla(\mathbf{p})$ is the volume of the hull we want to optimize as a function of \mathbf{p} estimated with Nautilus¹, ∇_l and ∇_u are the upper and lower bound of $\nabla(\mathbf{p})$. Finally, $R_t(\mathbf{p}, \cdot)$ is the total resistance as a function of \mathbf{p} and the velocity (computed at v^{High} and v^{Low}) estimated via the surrogate described in Section IV-A but where \mathbf{p} induces the geometry and, based on the geometry, Nautilus¹ estimates the quantities of Table 4 that together with the velocity are the actual inputs of the surrogate.

Problem (1) is a non-linear non-linearly constrained multi-objective optimization problem that is hard to optimize in practice.

The first step toward the solution of Problem (1) is to reformulate the problem as a single objective one. For this

TABLE 8. The parameter ranges for the 32 geometric design parameters of the values extracted from the DSYHS database. The parameter ranges are reported in mm.

Parameter	Range	Parameter	Range
X _{Bow}	[1650, 2200]	P _{Entr} ^{Y1}	[6, 43]
Y _{Bow}	[1, 62]	P _{Entr} ^{Z1}	[31, 85]
Z _{Bow}	[156, 355]	P _{Entr} ^{Y2}	[12, 71]
Z _{Max}	[196, 364]	P _{Entr} ^{Z2}	[60, 165]
X _{Entr}	[1600, 2050]	P _{Mid} ^{Y1}	[154, 349]
Y _{Entr}	[18, 97]	P _{Mid} ^{Z1}	[63, 121]
Z _{Entr}	[22, 179]	P _{Mid} ^{Y2}	[166, 368]
X _{Mid}	[800, 1050]	P _{Mid} ^{Z2}	[127, 244]
Y _{Mid}	[173, 371]	P _{Run} ^{Y1}	[72, 274]
Z _{Mid}	[−97, 4]	P _{Run} ^{Z1}	[28, 79]
X _{Run}	[−200, 50]	P _{Run} ^{Y2}	[128, 307]
Y _{Run}	[146, 327]	P _{Run} ^{Z2}	[64, 163]
Z _{Run}	[29, 180]	P _{Stern} ^{Y1}	[9, 248]
X _{Stern}	[150, 400]	P _{Stern} ^{Z1}	[4, 71]
Y _{Stern}	[26, 292]	P _{Stern} ^{Y2}	[17, 273]
Z _{Stern}	[38, 282]	P _{Stern} ^{Z2}	[7, 139]

purpose we will rely on a classical approach: replace the multiple objectives with a weighted sum of the different objectives (changing the sign in front to the objective so as to have all minimization or maximization) [72]

$$\begin{aligned} \min_{\mathbf{p}} \quad & \lambda \frac{R_t(\mathbf{p}, v^{\text{High}})}{\nabla(\mathbf{p})} + (1 - \lambda) \frac{R_t(\mathbf{p}, v^{\text{Low}})}{\nabla(\mathbf{p})}, \\ \text{s.t.} \quad & \nabla_l \leq \nabla(\mathbf{p}) \leq \nabla_u, \\ & \mathbf{p}_l(\delta) \leq \mathbf{p} \leq \mathbf{p}_u(\delta), \end{aligned} \quad (2)$$

where $\lambda \in [0, 1]$ defines the importance of the different objectives, i.e., for $\lambda \rightarrow 1$ we care more about $R_t(\mathbf{p}, v^{\text{High}})$ than $R_t(\mathbf{p}, v^{\text{Low}})$ and vice-versa for $\lambda \rightarrow 0$. Solving Problem (2) for different values of λ allows for the creation of the so-called Pareto frontier in a computationally efficient way [72].

Problem (2) is a non-linear non-linearly constrained optimization problem. In order to solve this problem different approaches can be exploited [73]. In the literature, there are a number of state-of-the-art algorithms available that are able to deal with this problem, e.g., gradient descent [74], swarm [75], and evolutionary [76]. A series of no-free-lunch theorems [77] ensure us that there is no way to choose a-priori the best optimization algorithms for a particular problem and the only option is to empirically test multiple approaches verifying which is actually the best one. As a consequence, to the best of the authors' knowledge and according to the literature on the subject [1], [16], [17], [19], [22], we opt for the Evolutionary Algorithm (EA) as it showed to be the best approach for these class of problems. In particular, we relied on an EA-based optimization framework built in MATLAB⁸ using the function `ga` which is a variant implementation of the NSGA-II [78], [79] Genetic Algorithm. Moreover,

⁸<https://mathworks.com/products/matlab.html>

TABLE 9. Parameters setting for the optimization algorithm exploited to solve Problem (2).

Algorithm	Matlab Function	Parameter	Value
EA	ga	PopulationSize	5000
		MaxGenerations	200
		CrossoverFraction	0.8
		EliteCount	1
		Multi start (manually implemented)	10

we customize the optimizer adding a multi-start approach, running the algorithm multiple times keeping the best solution found in the different starts. For the sake of repeatability, Table 9 reports the parameters' set that empirically produced the best results in the paper.

E. OPTIMIZATION FRAMEWORK PHYSICAL PLAUSIBILITY

In this section, we will present the proposed approach to demonstrate the physical plausibility of the solution (i.e., hull geometry) retrieved by solving Problem (1) through Problem (2) with different λ (see the previous section).

First, we need to better specify our definition of physical plausibility. In particular, in this work, we consider the ability of the optimizer to find non-degenerate geometries, namely geometries that in EFD will exhibit $R_t(\mathbf{p}, v^{\text{High}})$ and $R_t(\mathbf{p}, v^{\text{Low}})$ far away from the one suggested by the optimizer. Such geometry is then considered non-physically plausible. This outcome may happen for two main reasons, which are also connected

- the first one is because $R_t(\mathbf{p}, v^{\text{High}})$ and $R_t(\mathbf{p}, v^{\text{Low}})$ inserted in Problems (1) and (2) are not the real resistances but a surrogate characterized by no infinite precision and limited extrapolation abilities (this has been already tested in Section IV-B). As a consequence, during exploration, the EA can spot false minima induced by the imprecision and the extrapolation limitations of the surrogate model
- the second one is that the parameter space defined by $\mathbf{p}_l(\delta)$ and $\mathbf{p}_u(\delta)$, namely by δ is too large, requesting the optimizer to search within a parameter space that has more risk of imprecise extrapolation of the surrogate.

For this reason, analogously to what has been done for the surrogate in Section IV-B, we will test the geometries found by the optimizer with the Star CCM+⁶ package checking the deviation between the estimated $R_t(\mathbf{p}, v^{\text{High}})$ and $R_t(\mathbf{p}, v^{\text{Low}})$ and the one identified by the surrogate and then the optimizer. In the CFD simulation based on the Star CCM+⁶ package, we exploited the same setting described in Section IV-B.

V. EXPERIMENTAL RESULTS

In this section, we will report the results of applying the methodology described in Section IV to solve the problem faced in this work using the data described in Section III.

Specifically, we performed the following experiments

- in Section V-A we tested the quality of the surrogate model in the different extrapolating scenarios (LOVO, LOGO, and LOSO);
- in Section V-B we focused on the LOSO scenario, the most challenging and useful in practice, testing

the physical plausibility of the results against the CFD;

- in Section V-C we tested the quality of the optimization framework on a particular series of the DSYHS showing that we can improve the current geometries with new designs that we tested using CFD to verify their physical plausibility.

All experiments were performed with 2×Intel XEON E5-6248R 24C 3.0GHz CPUs and 192 GB of Memory.

A. SURROGATE MODELS VALIDATION IN THE EXTRAPOLATING SCENARIOS

In this section, we will report the performance of the surrogate models described in Section IV-A using the validation approaches described in Section IV-B in the different extrapolating scenarios. In particular, we will compare the results of the different algorithms employed to build the surrogate (RF, XGB, KRR, and ELM) on the different extrapolating scenarios (LOVO, LOGO, and LOSO) using different metrics. For the metrics, we measured the accuracy with both quantitative (MAE, MAPE, and PPMCC) and qualitative (scatter plot) measures and the computational requirements (Training Time and Test Time).

Table 10 reports for all algorithms employed to build the surrogate (RF, XGB, KRR, and ELM) and for all the different extrapolating scenarios (LOVO, LOGO, and LOSO) the different metrics employed to evaluate the performance (MAE, MAPE, PPMCC, Training Time, and Test Time). Figure 7, instead, reports the scatter plot for the best algorithm in each scenario (ELM for LOVO and KRR for LOGO and LOSO) where we considered just the Accuracy as a metric since the Test time differences are negligible.

From Table 10 and Figure 7 it is possible to observe that

- as the complexity of the extrapolation scenario increases (i.e., from LOVO to LOSO) the average accuracy of the models, across all of the algorithms, decreases;
- the ELM is the best performing algorithm for the LOVO scenario, while the KRR is the best performing algorithm for the LOGO and LOSO scenarios;
- despite the fact RF was demonstrated as the best algorithm overall in terms of Test Time, differences with the other methods are negligible for our application (well below fractions of milliseconds);
- final performance both in terms of accuracy (well below 1% of error) and Test Time (less than 10^{-5} [s]) even in the most challenging scenario (LOSO) make these surrogates perfect to be employed inside and automatic optimization framework (see Section IV-D). In fact, in order to reach this level of accuracy, usually a CFD simulation is required, but the same prediction takes around 1 hour with CFD.

B. SURROGATE VALIDATION AND PHYSICAL PLAUSIBILITY IN THE LOSO SCENARIO

In this section, we will deepen the analysis of the performance of the best algorithm identified in Section V-A for the

LOSO scenario (KRR), because, in practice, this is the most interesting scenario. In fact, in practice, what we want to do is to generate geometry for a new, previously unexplored series.

Let us start by validating the quality of the model on the different series and on the different geometries.

Table 11 reports, for the KRR in the LOSO scenario, the different metrics of accuracy (MAE, MAPE, and PPMCC) for each of the series. Instead, Figure 8 reports, for the KRR in the LOSO scenario, the scatter plot for each of the series.

From Table 11 and Figure 8 it is possible to observe that the surrogate performs better on some series than on others. This is due to several reasons

- the performance of the surrogate decreases at higher speeds as less experiments have been performed at higher speeds. As a matter of fact, for S_1 , S_2 , and S_4 , the poor performance is exhibited around Resistances in the range from 40–100N (see Figures 8a, 8b, and 8d);
- the poor performance for S_1 as the LOSO (Figure 8a) is related to the fact that there are significantly more geometries in this series than in any other (22 out of 47 geometries according to Table 3). Consequently, when we check S_1 in the LOSO scenario, we have very few geometries to learn our model.

Let us continue this section with the test of the physical plausibility of the surrogate.

First, we have to look in detail at the performance of the surrogate in each geometry of the series. Since reporting all the errors for all the geometries of all the series is not meaningful, we decided to report in Table 12, for the KRR in the LOSO scenario, the different metrics of accuracy (MAE, MAPE, and PPMCC) for the best (i.e., the one exhibiting the smallest error) and worst (i.e., the one exhibiting the most significant error) geometries in each of the series.

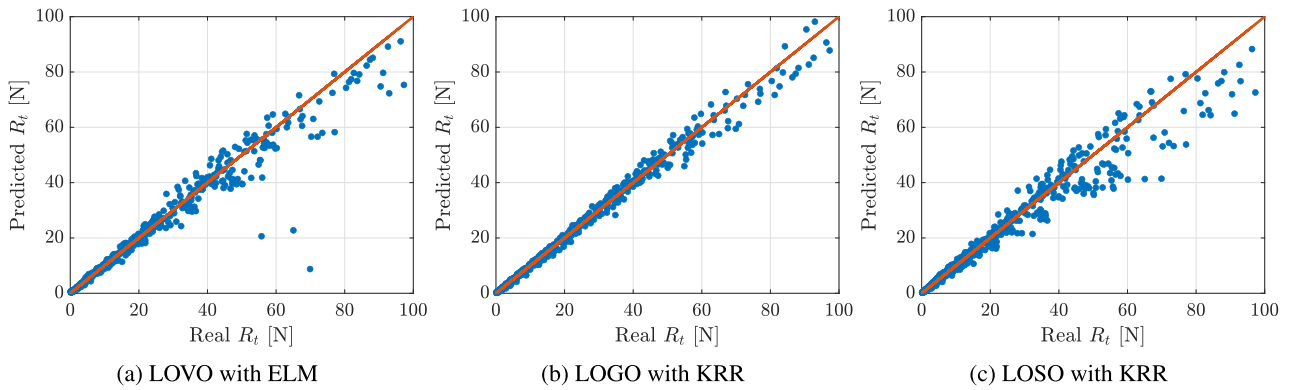
From Table 12, it is possible to observe how the gap between the best and the worst geometries can be significant. Nevertheless, large or small here is not a concept that we can define without having a baseline.

For this reason, in Figure 9, we reported, for the best and worst geometry in each of the series as in Table 12, the comparison between the EFD (the available data), the KRR surrogate (that we learned from the EFD in the LOSO scenario), and the CFD (using the StarCCM+⁶ package as described in Section IV-E). Note that, for the EFD just few points are actually known and we linearly interpolate between them. For the CFD, we have the same issue since making one prediction, as pointed out in Section V-A, takes a few hours. Instead, for the KRR, we can make prediction for a huge number of points since only fractions of milliseconds are needed (see Table 10).

From Table 12 and Figure 9, we can observe that the deviation of the KRR-based surrogate from the EFD is, in terms of magnitude, similar, when not better, than the one of the CFD even when we consider the geometry in which the surrogate performs worse. Moreover, the resistance behaviour as a function of the speed is quantitatively aligned with the expectations. In conclusion, the KRR-based

TABLE 10. Surrogate Models Validation in the Extrapolating Scenarios: metrics employed to evaluate the performance (MAE, MAPE, PPMCC, Training Time, and Test Time) for all algorithms employed to build the surrogate (RF, XGB, KRR, and ELM) and for all the different extrapolating scenarios (LOVO, LOGO, and LOSO).

Scenario	Surrogate	Accuracy			Time	
		MAE [N]	MAPE [%]	PPMCC [–]	Training 10^3 [s]	Test 10^{-6} [s]
LOVO	RF	2.56 ± 1.14	0.96 ± 0.62	0.89 ± 0.09	2.0 ± 0.1	2.2 ± 0.2
	XGB	2.35 ± 1.07	0.89 ± 0.63	0.96 ± 0.02	1.8 ± 0.1	5.0 ± 0.3
	KRR	4.75 ± 5.26	2.14 ± 2.77	0.85 ± 0.19	0.1 ± 0.1	9.6 ± 7.7
	ELM	2.28 ± 0.82	1.15 ± 0.45	0.94 ± 0.03	0.9 ± 0.0	5.2 ± 0.7
LOGO	RF	2.22 ± 1.39	0.12 ± 0.05	0.99 ± 0.00	1.9 ± 0.0	5.0 ± 0.3
	XGB	1.65 ± 0.97	0.12 ± 0.04	0.99 ± 0.00	1.6 ± 0.0	5.2 ± 0.4
	KRR	0.94 ± 0.21	0.48 ± 0.09	0.99 ± 0.01	0.5 ± 0.0	7.5 ± 0.5
	ELM	1.67 ± 0.35	0.63 ± 0.09	0.99 ± 0.00	5.9 ± 0.0	4.5 ± 0.5
LOSO	RF	2.84 ± 1.67	0.16 ± 0.06	0.99 ± 0.01	0.5 ± 0.0	1.2 ± 0.3
	XGB	2.33 ± 1.47	0.14 ± 0.08	0.99 ± 0.00	0.4 ± 0.0	4.8 ± 0.3
	KRR	1.83 ± 1.07	0.11 ± 0.03	0.99 ± 0.00	0.1 ± 0.0	7.6 ± 5.4
	ELM	3.11 ± 1.99	0.16 ± 0.07	0.96 ± 0.01	0.1 ± 0.0	1.5 ± 0.1

**FIGURE 7.** Surrogate Models Validation in the Extrapolating Scenarios: scatter plot for the best algorithm in each scenario (ELM for LOVO and KRR for LOGO and LOSO) considering just the Accuracy as a metric since the Test time differences are negligible (see Table 10).**TABLE 11.** Surrogate Validation in the LOSO Scenario: metrics of accuracy (MAE, MAPE, and PPMCC) for each of the series of the KRR (the best algorithm identified in Section V-A).

Series	Accuracy		
	MAE [N]	MAPE [%]	PPMCC [–]
S_1	3.47 ± 0.94	1.26 ± 0.29	0.99 ± 0.00
S_2	3.13 ± 1.85	0.61 ± 0.24	0.99 ± 0.02
S_3	0.62 ± 0.34	0.62 ± 0.35	1.00 ± 0.00
S_4	2.28 ± 0.49	0.88 ± 0.17	1.00 ± 0.00
S_6	0.86 ± 0.38	0.49 ± 0.11	1.00 ± 0.00
S_7	0.62 ± 0.18	0.37 ± 0.15	1.00 ± 0.00

TABLE 12. Surrogate Validation in the LOSO Scenario: different metrics of accuracy (MAE, MAPE, and PPMCC) for the best (i.e., the one exhibiting the smallest error) and worst (i.e., the one exhibiting the most significant error) geometries in each of the series for the KRR in the LOSO scenario.

Series	Best	Accuracy			Worst	Accuracy		
		MAE [N]	MAPE [%]	PPMCC		MAE [N]	MAPE [%]	PPMCC
S_1	G_7	1.50	0.60	0.99	G_6	10.47	1.83	0.99
S_2	G_4	0.34	0.29	0.99	G_5	5.31	1.00	0.99
S_3	G_1	0.11	0.14	0.99	G_7	1.47	1.52	0.99
S_4	G_9	0.77	0.86	0.99	G_5	3.63	1.42	0.99
S_6	G_3	0.58	0.37	0.99	G_1	0.98	0.53	0.99
S_7	G_1	0.51	0.27	0.99	G_2	0.72	0.44	0.99

surrogate performance and physical plausibility can be considered at the level of a state-of-the-art CFD-based model

at a fraction of its computational requirements: from hours to a fraction of milliseconds.

C. OPTIMIZATION FRAMEWORK VALIDATION AND PHYSICAL PLAUSIBILITY

At this point, we have empirically shown that the proposed parametrization and the surrogate are able to work well also in extrapolating scenarios matching the performance, in terms of accuracy and physical plausibility, of state-of-the-art CFD models at a fraction of their computational requirements. In this section, we will leverage this surrogate in the optimization framework proposed in Section IV-D, validating its performance by means of the approach described in Section IV-E.

For computational constraints (i.e., using the CFD too many times would result in months of simulations) in this section we limit the analysis to the optimization of a single series S_j . In order to have a realistic baseline (i.e., EFD data) we designed a specific experiment: we trained the surrogate with the EFD of all the series in the DSYHS except S_j simulating the need to design a vessel exactly in the missing series. In this way, the EFD data of S_j will function as a realistic baseline to compare with the results of our optimization. Note that, with this approach, we are actually using the surrogate as in the

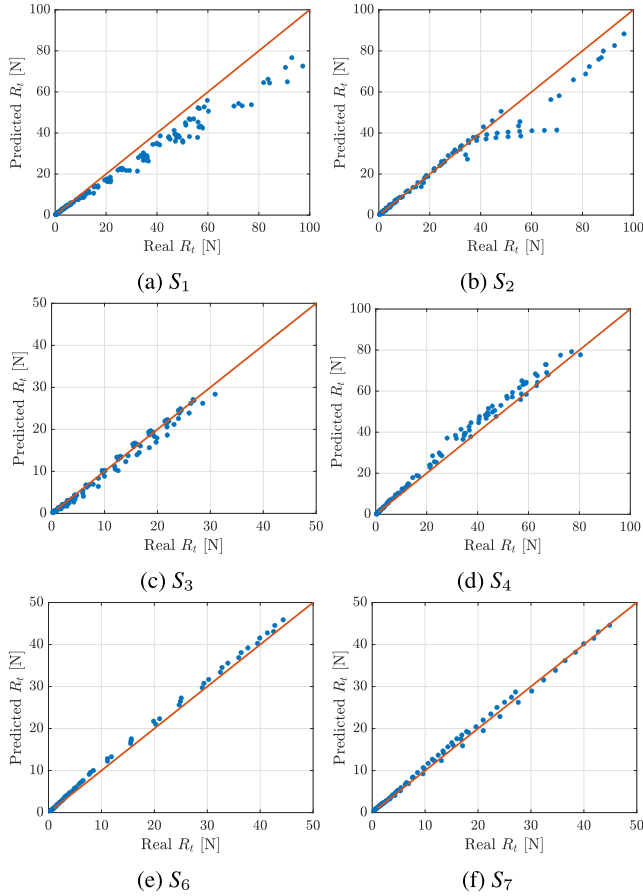


FIGURE 8. Surrogate Validation in the LOSO Scenario: scatter plot for the KRR (the best algorithm identified in Section V-A) for each of the series (see Table 10).

LOSO scenario. Then we solved the optimization problem of Section IV-D using this surrogate as resistance predictor for different values of $\lambda \in [0, 1]$ and with the parameter range induced by the S_j (see Section IV-D) computing the Pareto frontier of the geometries. The Pareto frontier of the geometries is then compared with the EFD data of the S_j (where we linearly interpolated between the available data). Moreover, for each one of the geometries on the Pareto we computed the resistance at high and low speed with the CFD.

We set $S_j = S_4$: this choice is based on Table 11 as this is the series that exhibits approximately the average performance of the surrogate in the LOSO scenario (i.e., it is not the most challenging nor the simplest series to predict but is an average to challenging one). For S_4 the $\mathbf{p}_l(\delta)$ and $\mathbf{p}_u(\delta)$ are reported in Table 8 while $\nabla_l = 19 \cdot 10^{-3} m^3$ and $\nabla_u = 48 \cdot 10^{-3} m^3$. We reported the results for different values of $\delta \in \{10, 20, 30\}\%$ and $\lambda \in \{0, 0.1, \dots, 1\}$ linearly interpolating between this value.

Figure 10 reports the Pareto frontier ($\frac{R_t(\mathbf{p}, v^{\text{High}})}{\nabla(\mathbf{p})}$ on the x-axis and $\frac{R_t(\mathbf{p}, v^{\text{Low}})}{\nabla(\mathbf{p})}$ on the y-axis) for different values of λ and δ together with the EFD data and the CFD validation as just described. Additionally, Figure 11 reports a comparison of the

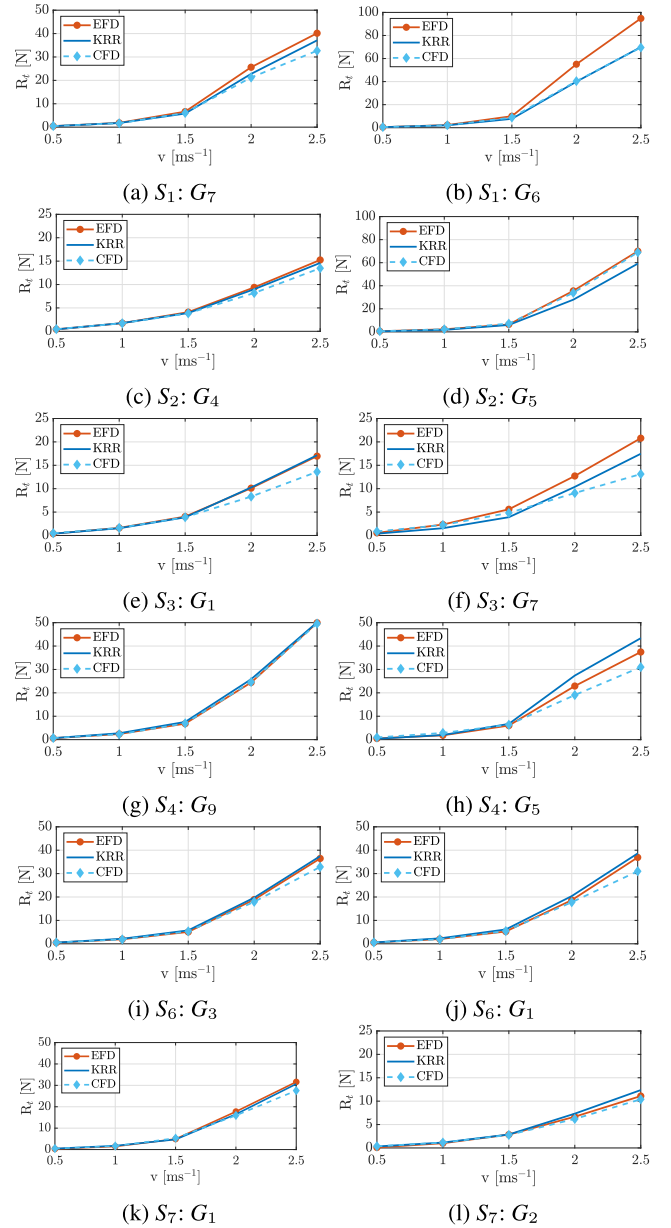


FIGURE 9. Surrogate Physical Plausibility in the LOSO Scenario: comparison between the EFD (the available data), the KRR surrogate (that we learned from the EFD in the LOSO scenario), and the CFD (using the StarCCM+⁶ package as described in Section IV-E) for the best and worst geometry in each of the series as in Table 12.

body plans⁹ for the baseline geometry belonging to S_4 and the optimized ones with $\lambda = 1$ and $\delta \in \{10, 20, 30\}\%$. Setting $\lambda = 1$ implies that we prefer to minimise the resistance at v^{High} , representing a typical velocity for high-speed operations where we should observe the most significant differences in optimal performance.

⁹The body plan is commonly used in naval architecture to display hull geometries and contains the set of transverse sections (the fore of the hull is on the left, and the aft on the right).

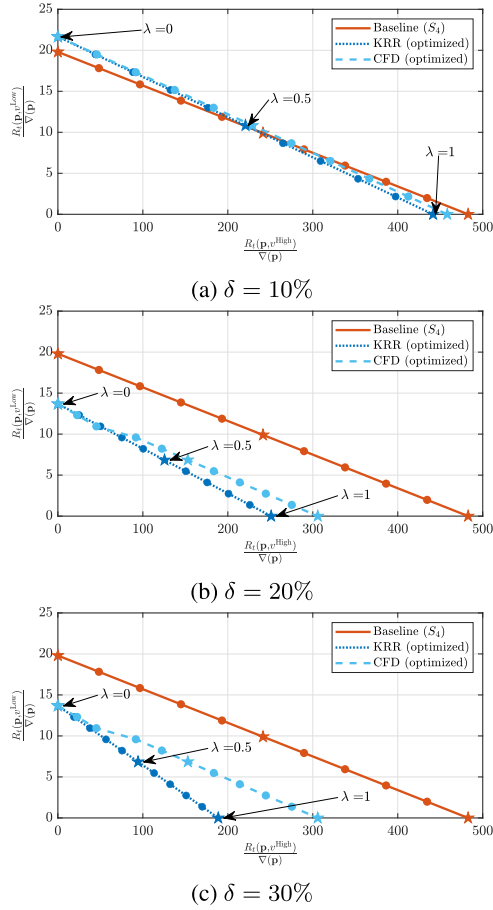


FIGURE 10. Optimization Framework Validation and Physical Plausibility: Pareto frontier ($\frac{R_t(p, v^{\text{High}})}{V(p)}$ on the x-axis and $\frac{R_t(p, v^{\text{Low}})}{V(p)}$ on the y-axis) for different values of λ and δ together with the EFD data and the CFD.

From Figures 10 and 11 we can observe that

- when δ is small (Figure 10a, $\delta = 10\%$) the optimization framework coupled with the surrogate is able to find geometries that match the performance of the one in S_4 without any a-priori knowledge of the geometries belonging to S_4 . Nevertheless, it is worth noting how the geometry found by the optimizer (Figure 11a), even if having a similar performance, is quite different. This is due to the fact that the optimization problem is surely simplified, not taking into account all the realistic constraints that impact the design of a hull geometry (e.g., stability and seakeeping);
- when δ is a bit larger (Figure 10b, $\delta = 20\%$) the surrogate is able to exceed remarkably, according to the surrogate, the performance of the S_4 geometry. However, this is a bit optimistic when checking the resistance at high speed: when using the CFD to estimate the resistance of the geometry found with the surrogate there is a reduction of this performance gain which remains still remarkable. Also in this case note that the differences in the geometries (Figure 11b) starts to enlarge;

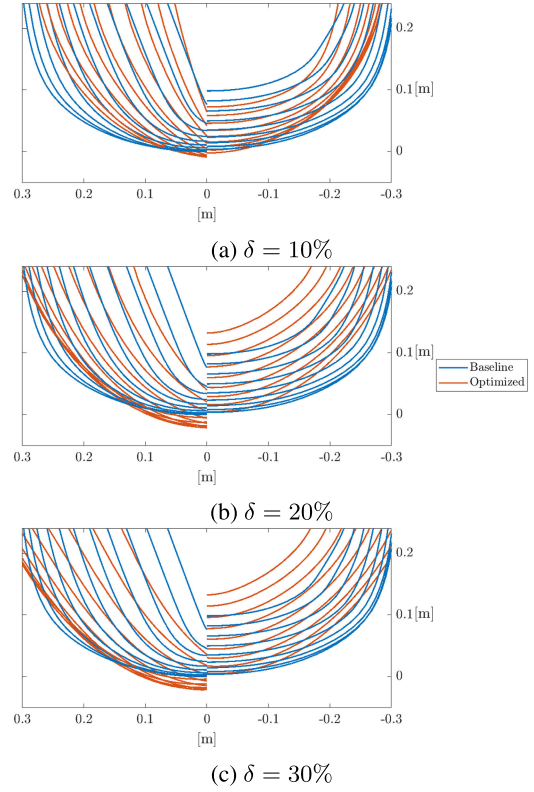


FIGURE 11. Optimization Framework Validation and Physical Plausibility: a comparison of the the body plans⁹ for the baseline geometry belonging to S_4 and the optimized ones with $\lambda = 1$ and $\delta \in \{10, 20, 30\}\%$.

- when we further increase δ (Figure 10c, $\delta = 30\%$) the surrogate can exceed even more, according to the surrogate, the performance of the S_4 geometry. However, this is just a numerical artefact when checking the resistance at high speed due to the extrapolation limits of the surrogate. In fact, when using the CFD to estimate the resistance of the geometry found with the surrogate, there is a reduction of this performance that brings us back to the gain found when δ was smaller. Note that in this case the geometry (Figure 11b) is quite similar to the case of $\delta = 20\%$ (Figure 11b).

Finally, for the sake of completeness, a qualitative indicator of the quality of the optimized geometry with $\lambda = 1$ (for the same reasons as before) and $\delta \in \{10, 20, 30\}\%$ is reported in Figure 12 which shows the wave profile at v^{High} of the original S_4 hull (top half) and the difference with the optimized parametric hulls (bottom half).

From Figure 12 we can observe that

- in all cases there is a noticeable difference between the original hulls and the optimized hulls;
- when $\delta = 10$ (Figure 12a) there is little significant difference (indicated by the lack of white color in the bottom half of the figure) which is expected due to the fact that the representation space is constrained around that of the original hull;
- when $\delta = 20$ or $\delta = 30$ (Figures 12b and 12c) there is a more significant difference between the original and optimized wave profiles (indicated by the presence of

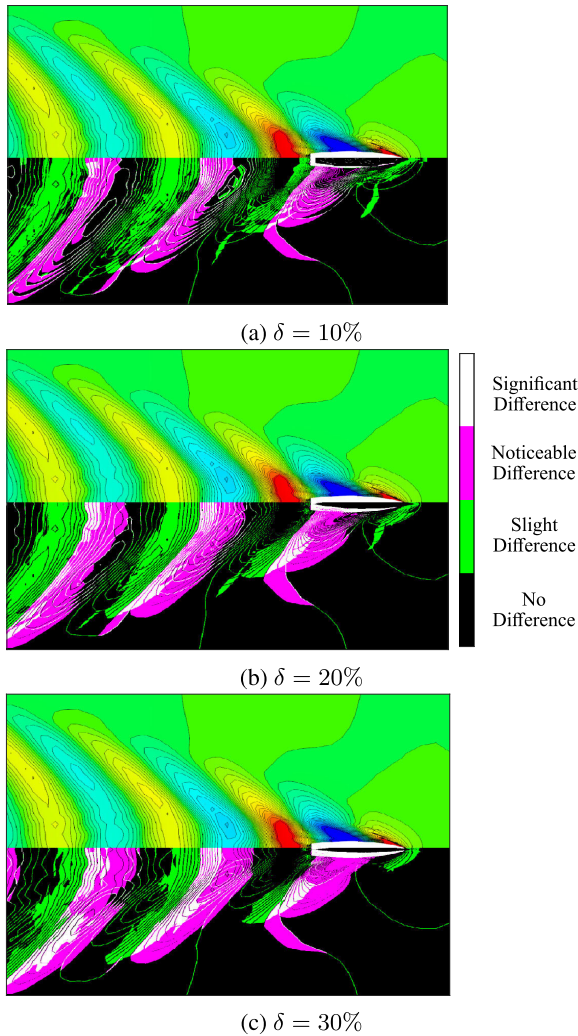


FIGURE 12. Optimization Framework Validation and Physical Plausibility: quality of the geometry generated in S_4 with $\lambda = 1$ and $\delta \in \{10, 20, 30\}\%$ using the wave profile at v^{High} of the baseline hull belonging to S_4 (top half) and the difference with the optimized one (bottom half).

white color in the bottom half of the figures) which is in agreement with the results of Figure 10.

VI. CONCLUSION

In this work, we tackled the problem of vessel hull resistance optimization, which is crucial for achieving optimal performance and reducing environmental impact. First, we reviewed the current approach in the literature that mostly relies on a mix of human experience and DDMs: human experts define, via parametrization and parameter ranges, a series of geometries; a surrogate of the relationship between these parameters and the resistance, based on data from EFD or CFD, is built to interpolate within the defined parameter ranges; finally, the optimal parameters are found by optimizing, with more or less human intervention, the surrogate and used to retrieve the optimal geometry. Several limitations of the existing approaches were identified, including the need for human intervention in geometry parametrization and optimization, extensive computational

efforts and costs, and limited ability to work beyond the specific settings. To overcome those limitations, and to the best of the authors' knowledge, for the first time in the literature, we proposed a parametrization able to accurately describe the entire DSYHS that was decoupled from the one needed to create the DDM. We showed that the DDM can be directly trained on the DSYHS EFD dataset, avoiding the need for new CFD or EFD customized for the specific problem, and match the performance of state-of-the-art CFD models even in extrapolating conditions (i.e., for geometries and parameter ranges beyond the boundaries used to construct the surrogate), with physical plausibility and minimal human intervention. Apart from the methodological contribution, we also validated our approach to developing DDMs on different and increasingly challenging extrapolating conditions with statistical methods using the DSYHS EFD dataset and for physical plausibility using state-of-the-art CFD models. We demonstrated the effectiveness of our proposal by showing that it is possible to optimize the hull resistance by exploring geometric parameters beyond the boundaries of the DSYHS.

ACKNOWLEDGMENT

The experimental portion of this work was conducted using DelftBlue supercomputer at Delft High-Performance Computing Center [80], which hosts 238 Compute nodes with a total of 476 Intel XEON E5-6248R 24C 3.0GHz CPUs and 192 GB of Memory per node.

REFERENCES

- [1] Y. Feng, Z. Chen, Y. Dai, F. Wang, J. Cai, and Z. Shen, "Multidisciplinary optimization of an offshore aquaculture vessel hull form based on the support vector regression surrogate model," *Ocean Eng.*, vol. 166, pp. 145–158, Oct. 2018.
- [2] S. Harries and C. Abt, "CAESES—The HOLISHIP platform for process integration and design optimization," in *A Holistic Approach to Ship Design*. Cham, Switzerland: Springer, 2019, pp. 247–293.
- [3] S. Skoupas, G. Zaraphonitis, and A. Papanikolaou, "Parametric design and optimisation of high-speed Ro-Ro Passenger ships," *Ocean Eng.*, vol. 189, Oct. 2019, Art. no. 106346.
- [4] V. N. Armstrong, "Vessel optimisation for low carbon shipping," *Ocean Eng.*, vol. 73, pp. 195–207, Nov. 2013.
- [5] N. Rehmatulla, J. Calleya, and T. Smith, "The implementation of technical energy efficiency and CO₂ emission reduction measures in shipping," *Ocean Eng.*, vol. 139, pp. 184–197, Jul. 2017.
- [6] J. J. Maisonneuve, S. Harries, J. Marzi, H. C. Raven, U. Viviani, and H. Piippo, "Towards optimal design of ship hull shapes," in *Proc. 8th Int. Mar. Design Conf.*, pp. 31–42, 2003.
- [7] S. N. Skinner and H. Zare-Behtash, "State-of-the-art in aerodynamic shape optimisation methods," *Appl. Soft Comput.*, vol. 62, pp. 933–962, Jan. 2018.
- [8] Z. Qiang, C. Hai-Chao, L. Zu-Yuan, F. Bai-Wei, Z. Cheng-Sheng, C. Xide, and W. Xiao, "Multi-stage design space reduction technology based on SOM and rough sets, and its application to hull form optimization," *Expert Syst. Appl.*, vol. 213, Mar. 2023, Art. no. 119229.
- [9] E. F. Campana, D. Peri, Y. Tahara, and F. Stern, "Shape optimization in ship hydrodynamics using computational fluid dynamics," *Comput. Methods Appl. Mech. Eng.*, vol. 196, nos. 1–3, pp. 634–651, Dec. 2006.
- [10] P. Temarel, W. Bai, A. Bruns, Q. Derbanne, D. Dessi, S. Dhavalikar, N. Fonseca, T. Fukasawa, X. Gu, A. Nestegård, A. Papanikolaou, J. Parunov, K. H. Song, and S. Wang, "Prediction of wave-induced loads on ships: Progress and challenges," *Ocean Eng.*, vol. 119, pp. 274–308, Jun. 2016.
- [11] J. Gerritsma, R. Onnink, and A. Versluis, "Geometry, resistance and stability of the Delft systematic yacht hull series," *Int. Shipbuilding Prog.*, vol. 28, no. 328, pp. 276–297, Dec. 1981.

- [12] J. A. Keuning, "Approximation of the hydrodynamic forces on a sailing yacht based on the 'Delft Systematic Yacht Hull Series,'" in *Proc. 15th Int. Symp. Yacht Design Yacht Construct.* Amsterdam, The Netherlands: WbMT, Nov. 1998, pp. 99–152.
- [13] J. A. Keuning and M. Katgert, "A bare hull resistance prediction method derived from the results of the delft systematic yacht hull series extended to higher speeds," in *Proc. Int. Conf. Innov. High Perform. Sailing Yachts.* Lorient, France, 2008, pp. 13–21.
- [14] F. Huang and C. Yang, "Hull form optimization of a cargo ship for reduced drag," *J. Hydrodynamics*, vol. 28, no. 2, pp. 173–183, Apr. 2016.
- [15] S. Zhang, T. Tezdogan, B. Zhang, L. Xu, and Y. Lai, "Hull form optimisation in waves based on CFD technique," *Ships Offshore Struct.*, vol. 13, no. 2, pp. 149–164, Feb. 2018.
- [16] J. Guerrero, A. Cominetti, J. Pralits, and D. Villa, "Surrogate-based optimization using an open-source framework: The bulbous bow shape optimization case," *Math. Comput. Appl.*, vol. 23, no. 4, p. 60, Oct. 2018.
- [17] A. Coppede, S. Gaggero, G. Vernengo, and D. Villa, "Hydrodynamic shape optimization by high fidelity CFD solver and Gaussian process based response surface method," *Appl. Ocean Res.*, vol. 90, Sep. 2019, Art. no. 101841.
- [18] K. Niklas and H. Prusko, "Full-scale CFD simulations for the determination of ship resistance as a rational, alternative method to towing tank experiments," *Ocean Eng.*, vol. 190, Oct. 2019, Art. no. 106435.
- [19] A. Miao and D. Wan, "Hull form optimization based on an NM+CFD integrated method for KCS," *Int. J. Comput. Methods*, vol. 17, no. 10, Dec. 2020, Art. no. 2050008.
- [20] P. Casalone, O. Dell'Edera, B. Fenu, G. Giorgi, S. A. Sirigu, and G. Mattiazzo, "Unsteady RANS CFD simulations of sailboat's hull and comparison with full-scale test," *J. Mar. Sci. Eng.*, vol. 8, no. 6, p. 394, May 2020.
- [21] S. Zhang, T. Tezdogan, B. Zhang, and L. Lin, "Research on the hull form optimization using the surrogate models," *Eng. Appl. Comput. Fluid Mech.*, vol. 15, no. 1, pp. 747–761, Jan. 2021.
- [22] X. Liu, W. Zhao, and D. Wan, "Multi-fidelity co-kriging surrogate model for ship hull form optimization," *Ocean Eng.*, vol. 243, Jan. 2022, Art. no. 110239.
- [23] Y. Tahara, J. Longo, and F. Stern, "Comparison of CFD and EFD for the series 60 $C_B=0.6$ in steady drift motion," *J. Mar. Sci. Technol.*, vol. 7, no. 1, pp. 17–30, Jun. 2002.
- [24] I. Biliotti, S. Brizzolara, M. Viviani, G. Vernengo, D. Ruscelli, M. Galliussi, D. Guadalupi, and A. Manfredini, "Automatic parametric hull form optimization of fast naval vessels," in *Proc. 11th Int. Conf. Fast Sea Transp. (FAST)*, Honolulu, HI, USA, 2011, pp. 294–301.
- [25] Y. Lin, Q. Yang, and G. Guan, "Automatic design optimization of SWATH applying CFD and RSM model," *Ocean Eng.*, vol. 172, pp. 146–154, Jan. 2019.
- [26] E. Lazarevska, "Comparison of different models for residuary resistance prediction," in *Proc. 9th EUROSIM Congr. Modelling Simulation (EUROSIM)*, 2016, pp. 511–517.
- [27] S. B. Šegota, N. Anđelić, J. Kudláček, and R. Čep, "Artificial neural network for predicting values of residuary resistance per unit weight of displacement," *J. Maritime Transp. Sci.*, vol. 57, no. 1, pp. 9–22, Dec. 2019.
- [28] S. B. Šegota, I. Lorencin, M. Šercer, and Z. Car, "Determining residuary resistance per unit weight of displacement with symbolic regression and gradient boosted tree algorithms," *Pomorstvo*, vol. 35, no. 2, pp. 287–296, Dec. 2021.
- [29] S. F. Fahrholz and J.-D. Caprace, "A machine learning approach to improve sailboat resistance prediction," *Ocean Eng.*, vol. 257, Aug. 2022, Art. no. 111642.
- [30] A. Zerbinati, A. Minelli, I. Ghazlane, and J. A. Désidéri, "Meta-model-assisted MGDA for multi-objective functional optimization," *Comput. Fluids*, vol. 102, pp. 116–130, Oct. 2014.
- [31] M. Mittendorf and A. D. Papanikolaou, "Hydrodynamic hull form optimization of fast catamarans using surrogate models," *Ship Technol. Res.*, vol. 68, no. 1, pp. 14–26, Jan. 2021.
- [32] Y. Lin, J. He, and K. Li, "Hull form design optimization of twin-skeg fishing vessel for minimum resistance based on surrogate model," *Adv. Eng. Softw.*, vol. 123, pp. 38–50, Sep. 2018.
- [33] F. Stern, J. Yang, Z. Wang, H. Sadat-Hosseini, M. Mousaviraad, S. Bhushan, and T. Xing, "Computational ship hydrodynamics: Nowadays and way forward," *Int. Shipbuilding Prog.*, vol. 60, pp. 3–105, Jul. 2013.
- [34] H. C. Raven, A. Van der Ploeg, A. R. Starke, and L. Eça, "Towards a CFD-based prediction of ship performance-progress in predicting full-scale resistance and scale effects," *Int. J. Maritime Eng.*, vol. 150, no. A4, 2008.
- [35] T. Hino, F. Stern, L. Larsson, M. Visonneau, N. Hirata, and J. Kim, *Numerical Ship Hydrodynamics: An Assessment of the Tokyo 2015 Workshop*, vol. 94. Springer, 2020.
- [36] J. Antony, *Design of Experiments for Engineers and Scientists*. Amsterdam, The Netherlands: Elsevier, 2014.
- [37] M. J. Anderson and P. J. Whitcomb, *RSM Simplified: Optimizing Processes Using Response Surface Methods for Design of Experiments*, 2nd ed. New York, NY, USA: Productivity Press, 2016.
- [38] Y. Tahara, D. Peri, E. F. Campana, and F. Stern, "Single- and multiobjective design optimization of a fast multihull ship: Numerical and experimental results," *J. Mar. Sci. Technol.*, vol. 16, no. 4, pp. 412–433, Dec. 2011.
- [39] H. Zhao, T. Icoz, Y. Jaluria, and D. Knight, "Application of data-driven design optimization methodology to a multi-objective design optimization problem," *J. Eng. Design*, vol. 18, no. 4, pp. 343–359, Aug. 2007.
- [40] R. S. Burachik, C. Y. Kaya, and M. M. Rizvi, "Algorithms for generating Pareto fronts of multi-objective integer and mixed-integer programming problems," *Eng. Optim.*, vol. 54, no. 8, pp. 1413–1425, Aug. 2022.
- [41] A. Serani and M. Diez, "Parametric model embedding," *Comput. Methods Appl. Mech. Eng.*, vol. 404, Feb. 2023, Art. no. 115776.
- [42] J. M. Walker, A. Coraddu, and L. Oneto, "A decoupled approach to AI-based design and optimization of the Delft systematic yacht hull series," in *Proc. 22nd Conf. Comput. IT Appl. Maritime Industries (COMPIT)*, Drübeck, Germany, May 2023, pp. 23–25.
- [43] M. Terziev, T. Tezdogan, and A. Incecik, "Scale effects and full-scale ship hydrodynamics: A review," *Ocean Eng.*, vol. 245, Feb. 2022, Art. no. 110496.
- [44] Delft University of Technology. (2021). *Delft Systematic Yacht Hull Series Database Website*. [Online]. Available: <http://dsyhs.tudelft.nl/>
- [45] J. de Baar, S. Roberts, R. Dwight, and B. Mallol, "Uncertainty quantification for a sailing yacht hull, using multi-fidelity kriging," *Comput. Fluids*, vol. 123, pp. 185–201, Dec. 2015.
- [46] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. Cambridge, U.K.: Cambridge Univ. Press, 2014.
- [47] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [48] S. P. Adam, S. A. N. Alexandropoulos, P. M. Pardalos, and M. N. Vrahatis, "No free lunch theorem: A review," in *Approximation and Optimization: Algorithms, Complexity and Applications*, 2019, pp. 57–82.
- [49] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do we need hundreds of classifiers to solve real world classification problems?" *J. Machine Learn. Res.*, vol. 15, no. 1, pp. 3133–3181, 2014.
- [50] M. Wainberg, B. Alipanahi, and B. J. Frey, "Are random forests truly the best classifiers?" *J. Machine Learn. Res.*, vol. 17, no. 1, pp. 3837–3841, 2016.
- [51] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [52] I. Orlandi, L. Oneto, and D. Anguita, "Random forests model selection," in *Proc. Eur. Symp. Artif. Neural Netw., Comput. Intell. Mach. Learn.*, 2016, pp. 441–446.
- [53] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. ASM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2016, pp. 785–794.
- [54] G.-B. Huang and L. Chen, "Convex incremental extreme learning machine," *Neurocomputing*, vol. 70, nos. 16–18, pp. 3056–3062, Oct. 2007.
- [55] G.-B. Huang, L. Chen, and C.-K. Siew, "Universal approximation using incremental constructive feedforward networks with random hidden nodes," *IEEE Trans. Neural Netw.*, vol. 17, no. 4, pp. 879–892, Jul. 2006.
- [56] C. M. Bishop, *Neural Networks for Pattern Recognition*. London, U.K.: Oxford Univ. Press, 1995.
- [57] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [58] S. S. Keerthi and C.-J. Lin, "Asymptotic behaviors of support vector machines with Gaussian kernel," *Neural Comput.*, vol. 15, no. 7, pp. 1667–1689, Jul. 2003.
- [59] M. Z. Naser and A. H. Alavi, "Error metrics and performance fitness indicators for artificial intelligence and machine learning in engineering and sciences," *Archit., Struct. Construction*, vol. 3, no. 4, pp. 499–517, Dec. 2023.
- [60] K. L. Sainani, "The value of scatter plots," *PM&R*, vol. 8, no. 12, pp. 1213–1217, Dec. 2016.
- [61] L. Oneto, *Model Selection and Error Estimation in a Nutshell*. Springer, 2020.

- [62] D. C. Wilcox, *Turbulence modeling for CFD*, vol. 2. La Canada, CA, USA: DCW Industries, 1998.
- [63] C. W. Hirt and B. D. Nichols, "Volume of fluid (VOF) method for the dynamics of free boundaries," *J. Comput. Phys.*, vol. 39, no. 1, pp. 201–225, Jan. 1981.
- [64] S. Harries, C. Abt, and K. Hochkirch, "Hydrodynamic modeling of sailing yachts," in *Proc. Chesapeake Sailing Yacht Symp.*, 2001, pp. 1–13.
- [65] F. Pérez and J. A. Suárez, "Quasi-developable -spline surfaces in ship hull design," *Computer-Aided Design*, vol. 39, no. 10, pp. 853–862, Oct. 2007.
- [66] F. L. Pérez, J. A. Clemente, J. A. Suárez, and J. M. González, "Parametric generation, modeling, and fairing of simple hull lines with the use of nonuniform rational B-spline surfaces," *J. Ship Res.*, vol. 52, no. 1, pp. 1–15, Mar. 2008.
- [67] F. Pérez-Arribas, "Parametric generation of planing hulls," *Ocean Eng.*, vol. 81, pp. 89–104, May 2014.
- [68] S. Khan, E. Gunpinar, and K. M. Dogan, "A novel design framework for generation and parametric modification of yacht hull surfaces," *Ocean Eng.*, vol. 136, pp. 243–259, May 2017.
- [69] S. Ö. Felek, "Parametric sailing yacht exterior and interior design," *Tasarim Kuram*, vol. 16, no. 29, pp. 1–15, 2020.
- [70] A. Mancuso, A. Saporito, and D. Tumino, "Parametric hull design with rational Bézier curves," in *Advances on Mechanics, Design Engineering and Manufacturing III: Proceedings of the International Joint Conference on Mechanics, Design Engineering & Advanced Manufacturing, JCM 2020, June 2–4, 2020*. Springer, 2021, pp. 221–227.
- [71] S. M. Shamsuddin, M. A. Ahmed, and Y. Smian, "NURBS skinning surface for ship hull design based on new parameterization method," *Int. J. Adv. Manuf. Technol.*, vol. 28, nos. 9–10, pp. 936–941, Jul. 2006.
- [72] M. T. M. Emmerich and A. H. Deutz, "A tutorial on multiobjective optimization: Fundamentals and evolutionary methods," *Natural Comput.*, vol. 17, no. 3, pp. 585–609, Sep. 2018.
- [73] M. J. Kochenderfer and T. A. Wheeler, *Algorithms for Decision Making*. Cambridge, MA, USA: MIT Press, 2022.
- [74] S. Ruder, "An overview of gradient descent optimization algorithms," 2016, *arXiv:1609.04747*.
- [75] R. Poli, J. Kennedy, and T. Blackwell, "Particle swarm optimization: An overview," *Swarm Intell.*, vol. 1, pp. 33–57, Aug. 2007.
- [76] P. A. Vikhar, "Evolutionary algorithms: A critical review and its future prospects," in *Proc. Int. Conf. Global Trends Signal Process., Inf. Comput. Commun. (ICGTSPICC)*, Dec. 2016, pp. 261–265.
- [77] D. H. Wolpert and W. G. Macready, "No free lunch theorems for optimization," *IEEE Trans. Evol. Comput.*, vol. 1, no. 1, pp. 67–82, Apr. 1997.
- [78] J. McCall, "Genetic algorithms for modelling and optimisation," *J. Comput. Appl. Math.*, vol. 184, no. 1, pp. 205–222, Dec. 2005.
- [79] K. Deb, *Multiobjective Optimization Using Evolutionary Algorithms*. Hoboken, NJ, USA: Wiley, 2001.
- [80] Delft High Performance Computing Centre (DHPC). (2022). *DelftBlue Supercomputer (Phase 1)*. [Online]. Available: <https://www.tudelft.nl/dhpc/ark>



JAKE M. WALKER (Member, IEEE) was born in Edinburgh, U.K., in 1997. He received the M.Eng. degree (Hons.) in mechanical engineering from the University of Strathclyde, U.K., in 2020. He is currently pursuing the Ph.D. with the Department of Maritime and Transport Technology, Delft University of Technology, The Netherlands, with his primary focus on developing artificial intelligence-based solutions for the design and optimization of hull forms. In October 2020, he began a research position with the Naval Architecture and Marine Engineering Department, University of Strathclyde, where he transitioned his skills into the maritime domain.



ANDREA CORADDU (Member, IEEE) was born in Pietrasanta, Italy, in 1979. He received the Laurea degree in naval architecture and marine engineering from the University of Genoa, Italy, in 2006, and the Ph.D. degree from the School of Fluid and Solid Mechanics, University of Genoa, in 2012. His Ph.D. dissertation was titled, "Modeling and Control of Naval Electric Propulsion Plants." He was an Associate Professor with the Department of Naval Architecture, Ocean and Marine Engineering, University of Strathclyde, from October 2020 to August 2021. Currently, he is an Associate Professor of intelligent and sustainable energy systems with the Maritime and Transport Technology Department, Delft University of Technology, Delft, The Netherlands. His relevant professional and academic experiences, include working as an Assistant Professor with the University of Strathclyde, a Research Associate with the School of Marine Science and Technology, Newcastle University, and a Research Engineer as part of the DAMEN Research and Development Department, Singapore. He is also a Postdoctoral Research Fellow with the University of Genoa. He has been involved in a number of successful grant applications from research councils, industry, and international governmental agencies focusing on the design, integration, and control of complex marine energy and power management systems enabling the development of next-generation complex and multi-function vessels that can meet the pertinent social challenges regarding the environmental impact of human-related activities.



LUCA ONETO (Senior Member, IEEE) was born in Rapallo, Italy, in 1986. He received the B.Sc. and M.Sc. degrees in electronic engineering from the University of Genoa, Italy, in 2008 and 2010, respectively, and the Ph.D. degree from the School of Sciences and Technologies for Knowledge and Information Retrieval, University of Genoa, in 2014, with the thesis "Learning Based On Empirical Data." He was an Assistant Professor of computer engineering with the University of Genoa, from 2016 to 2019. In 2017, he obtained Italian National Scientific Qualification for the role of an Associate Professor of computer engineering and in 2018, he obtained the one in computer science. In 2018, he was the Co-Funder of the spin-off ZenaByte s.r.l. In 2019, he obtained Italian National Scientific Qualification for the role of a Full Professor of computer science and computer engineering. In 2019, he became an Associate Professor of computer science with the University of Pisa. He is currently an Associate Professor of computer engineering with the University of Genoa. He has been involved in several H2020 projects (S2RJU, ICT, and DS). His first main topic of research is the statistical learning theory with a particular focus on the theoretical aspects of the problems of (semi) supervised model selection and error estimation. His second main topic of research is data science with particular reference to the problem of trustworthy AI and the solution to real world problems by exploiting and improving the most recent learning algorithms and theoretical results in the fields of machine learning and data mining. He has been awarded with the Amazon AWS Machine Learning and Somalvico (Best Italian Young AI Researcher) Awards.

...