

Kevin Andrés Huilca Aparicio

# Structural NanoFingerprint Generation

MASTER'S DEGREE FINAL PROJECT

Directed by Dr. Francesc Serratosa

Master's degree in Computer Security Engineering and Artificial Intelligence



UNIVERSITAT ROVIRA I VIRGILI

Tarragona

2023

## **Abstract**

The goal of this project is to focus on the development of a publicly accessible website application to aid in the generation of NanoFingerprints, which will later be used to train models for toxicity prediction of different nanocompounds. In this paper a description of the development of the web portal has been included as well as the different functionalities that a user can interact with. Additionally, the NanoFingerprint generated have been used to train models for toxicity prediction by means of Linear Regression. Models that have shown better performance have been added to the website as means of providing additional tools for aiding in the investigation of nanocompounds. This project has been developed under the SbD4Nano funded project together with Universitat Rovira I Virgili.

## **Resum**

La finalitat d'aquest projecte és treballar en el desenvolupament d'un portal web de lliure accés al públic que permeti generar NanoFingerprints, que seran emprats per entrenar models per predir la toxicitat de diferents nanocomponents. En aquest document s'ha descrit el desenvolupament del portal web tenint en compte les funcionalitats disponibles amb les quals un usuari pot interactuar. A més, els NanoFingerprints generats s'han fet servir per entrenar models de predicción de toxicitats mitjançant Regressions Lineals. Els models que han mostrat millors resultats han estat afegit al portal web per tal de proveir d'eines addicional per ajudar amb la investigació de nanocomponents. Aquest projecte ha estat desenvolupat conjuntament amb la Universitat Rovira i Virgili sota el projecte finançat de SbD4Nano.

## **Resumen**

La finalidad de este proyecto ha sido el desarrollo de un portal web de libre acceso al público que permitiese generar NanoFingerprints, que luego se utilizan para entrenar modelos de predicción de toxicidad de diferentes nanocomponentes. En este documento se ha descrito el desarrollo de la página web, teniendo en cuenta las funcionalidades disponibles que son accesibles a nivel usuario. Además, los NanoFingerprints generados se han utilizado para entrenar dichos modelos de predicción mediante el uso de Regresiones Lineales. Los modelos que han dado mejores resultados han sido añadidos en la página web para proveer de herramientas adicionales para ayudar en la investigación de nanocomponentes. Este Proyecto ha sido desarrollado juntamente con la Universitat Rovira I Virgili bajo el proyecto financiado de SbD4Nano.

## *Acknowledgements*

I would like to thank Dr. Francesc Serratosa for his support when working through the different milestones of this project. Additionally, I am grateful to have had the opportunity to take part in a European project such as Sb4DNano. Finally, I would like to thank my family and friends for their patience and understanding at all times.

# Table of Contents

<b>Abstract</b> .....	2
<b>Resum</b> .....	2
<b>Resumen</b> .....	2
<b>Acknowledgements</b> .....	3
<b>Table of Figures</b> .....	5
<b>1. Introduction</b> .....	7
<b>1.1. Structure of thesis</b> .....	8
<b>2. Related work</b> .....	8
<b>2.1. Background</b> .....	8
<b>2.2. Nanofingerprint structure</b> .....	9
<b>3. Methodology</b> .....	12
<b>3.1. Atena Web Portal</b> .....	12
<b>3.1.1. Characteristics</b> .....	13
<b>3.1.2. Development environment</b> .....	14
<b>3.2. Webpage sections</b> .....	15
<b>3.2.1. NanoFingerprint</b> .....	15
<b>3.2.2. Toxicity Prediction</b> .....	19
<b>3.2.3. NanoFingerprint Toxicity Prediction</b> .....	23
<b>3.2.4. Subcomponent Search</b> .....	25
<b>3.2.5. Examples</b> .....	26
<b>3.2.6. Data Input Format</b> .....	27
<b>3.3. Nanofingerprint Generation</b> .....	27
<b>3.3. Dataset</b> .....	29
<b>3.4. Metrics</b> .....	30
<b>3.3.1. CCC</b> .....	30
<b>3.3.2. RMSE</b> .....	30
<b>3.3.3. MAE</b> .....	30
<b>3.3.4. R<sup>2</sup></b> .....	30
<b>4. Experimental Validation</b> .....	31
<b>4.1. Environment</b> .....	31
<b>4.2. Web Portal Development</b> .....	31
<b>4.3. Data Preprocessing</b> .....	32
<b>4.3.1. Papa – TiO<sub>2</sub> dataset</b> .....	33
<b>4.3.2. Papa – ZnO dataset</b> .....	36

4.3.3. Papa – TiO <sub>2</sub> +ZnO dataset .....	36
4.3.4. Anantha dataset .....	37
4.4. Model results.....	37
4.4.1. Papa – TiO <sub>2</sub> model .....	37
4.4.2. Papa – ZnO model .....	40
4.4.3. Papa – TiO <sub>2</sub> +ZnO model.....	42
4.4.4. Anantha model.....	44
6. Conclusions & Future Work .....	47
7. References.....	48
<b>Appendix A .....</b>	<b>50</b>
Access to Atena and Aura .....	50
Web Portal: GenVector function .....	51
Code model analysis .....	51
Models Generated .....	52
<b>Appendix B .....</b>	<b>53</b>
RFE Results .....	53
TiO <sub>2</sub> .....	53
ZnO .....	55
TiO <sub>2</sub> +ZnO.....	61

## Table of Figures

Figure 1. Project diagram.....	7
Figure 2. NanoFingerprint Sections 1 & 2 .....	10
Figure 3. NanoFingerprint Sections 3 & 4 .....	11
Figure 4. Nanocompound representation example.....	12
Figure 5. Project structure .....	13
Figure 6. Atena web portal homepage .....	15
Figure 7. NanoFingerprint generation main page.....	15
Figure 8. Error catch example .....	16
Figure 9. Results availability example .....	16
Figure 10. Verbose NanoFingerprint example .....	17
Figure 11. Nanofingerprint example.....	17
Figure 12. XYZ file with shell-only atoms .....	18
Figure 13. XYZ file with labels indicating shell atoms. ....	18
Figure 14. Toxicity Prediction.....	19
Figure 15. LDH (TiO <sub>2</sub> + ZnO) toxicity calculation.....	19
Figure 16. LDH (TiO <sub>2</sub> ) toxicity calculation .....	19
Figure 17. LDH (ZnO) toxicity calculation .....	20
Figure 18. Example of incorrect input parameters submission .....	20

Figure 19. Example of output for LDH (TiO2).....	20
Figure 20. Anantha toxicity prediction.....	21
Figure 21. Anantha prediction example.....	21
Figure 22. Toxicity prediction using equation in [6].....	22
Figure 23. Toxicity prediction using equation in [5].....	22
Figure 24. Toxicity prediction for (TiO <sub>2</sub> +ZnO) model.....	22
Figure 25. Toxicity models for TiO <sub>2</sub> and ZnO .....	23
Figure 26. NanoFingerprint toxicity prediction options.....	23
Figure 27. TiO <sub>2</sub> + ZnO + NanoFingerprint .....	24
Figure 28. TiO <sub>2</sub> + NanoFingerprint.....	24
Figure 29. ZnO + NanoFingerprint.....	24
Figure 30. Anantha equation + NanoFingerprint .....	25
Figure 31. Subcomponent search section.....	25
Figure 32. Substructure search result example .....	26
Figure 33. Directory tree for the Modelling and Subcomponent examples .....	27
Figure 34. Expected input files details. ....	27
Figure 35. Form input variables parsing .....	28
Figure 36. Go routine call to CoVectorCalc .....	28
Figure 37. TiO <sub>2</sub> dataset example (Part 1) .....	33
Figure 38. TiO <sub>2</sub> dataset example (Part 2) .....	33
Figure 39. TiO <sub>2</sub> joined sections dataset.....	34
Figure 40. TiO <sub>2</sub> _30 Verbose NanoFingerprint .....	34
Figure 41. TiO <sub>2</sub> training dataset example by joining sections .....	35
Figure 42. TiO <sub>2</sub> dataset after preprocessing example .....	35
Figure 43. TiO <sub>2</sub> training dataset example (Part 1) .....	36
Figure 44. TiO <sub>2</sub> training dataset example (Part 2) .....	36
Figure 45. Model metrics using TiO <sub>2</sub> dataset joining sections.....	37
Figure 46.TiO <sub>2</sub> modelling results using proposed features by authors .....	38
Figure 47.TiO <sub>2</sub> modelling results without features with std=0.....	38
Figure 48.TiO <sub>2</sub> modelling results with proposed features in article.....	38
Figure 49. RFE applied to TiO <sub>2</sub> dataset .....	39
Figure 50. ZnO modelling results using joined sections .....	40
Figure 51. ZnO modelling results using proposed features in article.....	40
Figure 52. ZnO modelling results without features with std = 0.....	41
Figure 53. ZnO modelling results using proposed features in article.....	41
Figure 54. RFE applied to ZnO dataset .....	41
Figure 55. TiO <sub>2</sub> +ZnO modelling results when joining sections .....	42
Figure 56. TiO <sub>2</sub> +ZnO modelling results using proposed features in article .....	42
Figure 57. TiO <sub>2</sub> +ZnO modelling results without features with std = 0 .....	43
Figure 58. TiO <sub>2</sub> +ZnO modelling results when using proposed features in article .....	43
Figure 59. RFE applied to TiO <sub>2</sub> +ZnO dataset .....	44
Figure 60. Anantha modelling results .....	44
Figure 61. Compiled results for TiO <sub>2</sub> .....	45
Figure 62. Compiled results for ZnO .....	46
Figure 63. Compiled results for TiO <sub>2</sub> +ZnO .....	46

## 1. Introduction

The following project has been developed under the SbD4Nano<sup>1</sup> funded project together with Universitat Rovira I Virgili. The presented work is divided in 2 main sections, as part of the study presented in [1], which aims to provide alternative toxicity prediction models by the analysis of nanocompounds as attributed graphs. These sections are:

- Development of a public website application
- Presentation of alternative toxicity prediction models using NanoFingerprints

A representation of the work presented in this thesis is shown below:

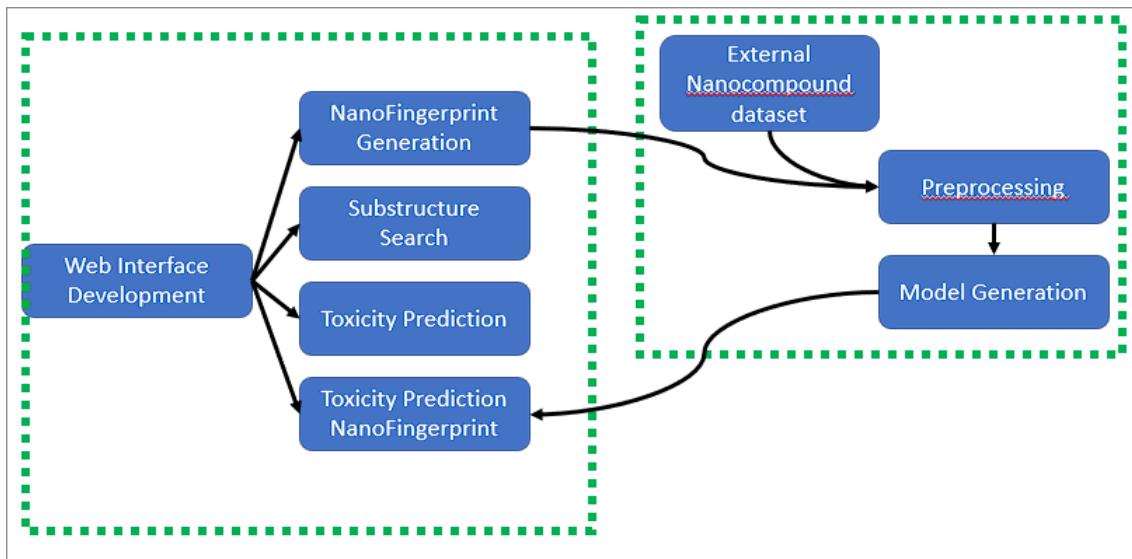


Figure 1. Project diagram

In terms of the web portal, the requirement is to develop the publicly accessible framework in charge of generating NanoFingerprints, ensuring all sections present in the definition of the NanoFingerprint are included considering the hardware constraints linked to the server hosting the site. NanoFingerprint is defined as a vector generated as a result of the analysis of a nanocompound considering it a graph.

Second part of the project consists of using the generated NanoFingerprints to define new models for toxicity prediction. Website will host known models and will include at the same time the proposed ones using NanoFingerprints, offering an alternative to the currently present in the literature.

---

<sup>1</sup> <https://www.sbd4nano.eu/project-overview>

### 1.1. Structure of thesis

This document is structured as follows: Section 2 corresponds to the presentation of the algorithms that are going to be analysed and compared to when proposing the toxicity models including the generated NanoFingerprints. The definition of the NanoFingerprint is also included, as well as an explanation of the different elements that form the array. A mention on the models that are going to be part of the case study are also described. Section 3 focuses on the presentation of the web portal as well as the different elements that form it, both at a backend and a frontend level. Additionally, a description of the datasets obtained joining both samples included in the literature as well as the NanoFingerprints generated is included. Furthermore, a definition of the metrics used to validate the models are shared. In Section 4, the preprocessing steps as well as the results associated with the models generated are included. Finally, Section 5 a discussion on the conclusions extracted from the project as well as possible future work development is presented.

## 2. Related work

### 2.1. Background

Multiple research involving the investigation of nanoparticles effect on living organisms has been presented, to have a better understanding not only of the implications they may have in nature but also considering the increase in use at both a commercial and industrial level [2]. As part of this project, 3 different approaches to analyse toxicity levels of a given nanocompound are explored and analysed.

The approach presented in [3] provides multiple MLR analyses in order to predict the impact  $TiO_2$  and  $ZnO$  have on cells when inducing the release of lactate dehydrogenase (LDH) which is directly related to damage on tissue membrane cells when in direct contact. Three different MLR models are described, included below:

$$LDH(TiO_2 + ZnO) = 0.66 + 0.005 * X_0 - 5 * X_2 + 0.003 * X_4 \quad (1)$$

$$LDH(TiO_2) = 0.599 + 0.004 * X_0 + 0.003 * X_4 \quad (2)$$

$$LDH(ZnO) = 1.041 + 0.001 * X_1 + 0.001 * X_2 + 0.001 * X_4 \quad (3)$$

$X_0$  refers to the size in nanometres of the nanocompound,  $X_1$  is listed as the size of the same nanocompound in water,  $X_2$  corresponds to the size of the nanocompound in Phosphate buffered saline (PBS), and  $X_4$  is defined as the concentration of the sample, measured in mg/l.

The toxicity analysis presented in [4] is based on multiple nanocompounds, in this case other qualities of the samples are considered. No prediction of LDH is considered, instead it is stated that if the output of the MLR is >0.5, the nanocomponent can be considered as TOXIC, otherwise it will be considered NOT TOXIC. Below is included the Linear regression model, which was included in NanoTox<sup>2</sup> nanotoxicology pipeline:

$$\begin{aligned} \text{Toxicity} = & 29.964 - \text{Size} * 4.39 + \text{HydroSize} * 1.564 - \text{SurfaceCharge} * \\ & 1.7914 - \text{SurfaceArea} * 4.754 - \text{Ec} * 12.662 + \text{Time} * 1.0105 + \text{Dose} * 6.259 + \\ & \text{Eneg} * 7.983 - \text{Noxygen} * 10.557 \end{aligned} \quad (4)$$

In both [5] [6], study is focused on the toxicity level for metal-oxide samples and their impact on human cells. Models proposed differ in both as a different enthalpy is used, below additional details regarding both models:

$$Puzyn - \log\left(\frac{1}{EC_{50}}\right) = 2.59 + -50 * \text{EnthalpyGaseousCation} \quad (5)$$

$$Gajewicz - \log\left(\frac{1}{EC_{50}}\right) = 2.47 + 0.24 * \text{EnthalpyMetalOxide} + 0.39 * X^c \quad (6)$$

$X^c$  is defined as the Mulliken's electronegativity (in eV). All these parameters are included in the datasets extracted from the included articles.

## 2.2. Nanofingerprint structure

Presented in [1], a NanoFingerprint is a vector representation generated as a result of the analysis of a nanocompound. This nanocomponent can be understood as an attributed graph, by considering atoms as nodes and chemical bonds as edges between nodes. Given the information that can be attached to these nodes as well as the mathematical characteristics linked to graphs, a substructure graph is defined only using the shell of the nanoparticles, which is known to hold most of the information linked to the toxicity of the compound [7].

The NanoFingerprint array can be divided in subsections, based on the relationship between nodes and the group of subgraphs that are generated. These subgraphs are defined as follows:

- $O(x)$ : number of oxygen atoms that have  $x$  bonds.
- $M(x)$ : number of metal atom that have  $x$  bonds.

---

<sup>2</sup> <https://github.com/NanoTox>

- $O(x, y)$ : subgraph with a central oxygen atom connected to  $x$  oxygen atoms and  $y$  metal atoms.
- $M(x, y)$ : subgraph with a central metal atom connected to  $x$  oxygen atoms and  $y$  metal atoms.
- $O(x, y)-O(x', y')$ : subgraph composed of  $O(x, y)$  and  $O(x', y')$  with central oxygen atoms connected.
- $M(x, y)-M(x', y')$ : subgraph composed of  $M(x, y)$  and  $M(x', y')$  with central metal atoms connected.
- $O(x, y)-M(x', y')$ : subgraph composed of  $O(x, y)$  and  $M(x', y')$  with central oxygen atom connected to central metal atom.

Considering these definitions, subsections are presented below. Note that index of the different elements have been added as well for better understanding of the size of the NanoFingerprint:

## Section 1

1. Shell thickness in Angstrom
2. Maximum number of bond per atom: **MAX**
3. Size in Angstrom
4. Atomic number of the metal
5. Number of oxygen atoms
6. Number of metal atoms
7. Number of O(1)
- ...

## Section 2

- MAX+6. Number of O(**MAX**)
- MAX+7. Number of M(**1**)
- ...
- 2MAX+6. Number of M(**MAX**)

Figure 2. NanoFingerprint Sections 1 & 2

<b>Section 3</b>	2MAX+7. Number of O(0,1) ... (MAX+1) <sup>2</sup> +2MAX+6. Number of O(MAX, MAX) (MAX+1) <sup>2</sup> +2MAX+7. Number of M(0, 1) ... 2(MAX+1) <sup>2</sup> +2MAX+6. Number of M(MAX, MAX)
<b>Section 4</b>	2(MAX+1) <sup>2</sup> +2MAX+7. Number of O(0,1)-O(0,1) ... (MAX+1) <sup>4</sup> +2(MAX+1) <sup>2</sup> +2MAX+6. Number of O(MAX,MAX)-O(MAX,MAX) (MAX+1) <sup>4</sup> +2(MAX+1) <sup>2</sup> +2MAX+7. Number of M(0,1)-M(0,1) ... 2(MAX+1) <sup>4</sup> +2(MAX+1) <sup>2</sup> +2MAX+6. Number of M(MAX,MAX)-M(MAX,MAX) 2(MAX+1) <sup>4</sup> +2(MAX+1) <sup>2</sup> +2MAX+7. Number of O(0,1)-M(0,1) ... 3(MAX+1) <sup>4</sup> +2(MAX+1) <sup>2</sup> +2MAX+6. Number of O(MAX,MAX)-M(MAX,MAX)

Figure 3. NanoFingerprint Sections 3 & 4

Considering the sections above, it can be stated that the maximum length of a NanoFingerprint is linked to MAX value, which is the maximum number of bonds to be analysed. As such, the length of the NanoFingerprint is computed considering the following equation:

$$\text{NanoFingerprint length} = 3 * (MAX + 1)^4 + 2 * (MAX + 1)^2 + 2 * MAX + 6 \quad (1)$$

MAX Value	NanoFingerprint Length
4	1939
5	3976
6	7319
7	12436
8	19869
9	30224
10	44191

Table 1. Impact on NanoFingerprint length given a MAX value

For A visual representation of a nanocompound has been added below to aid in the understanding of the different substructures that are analysed and extracted to form the NanoFingerprint:

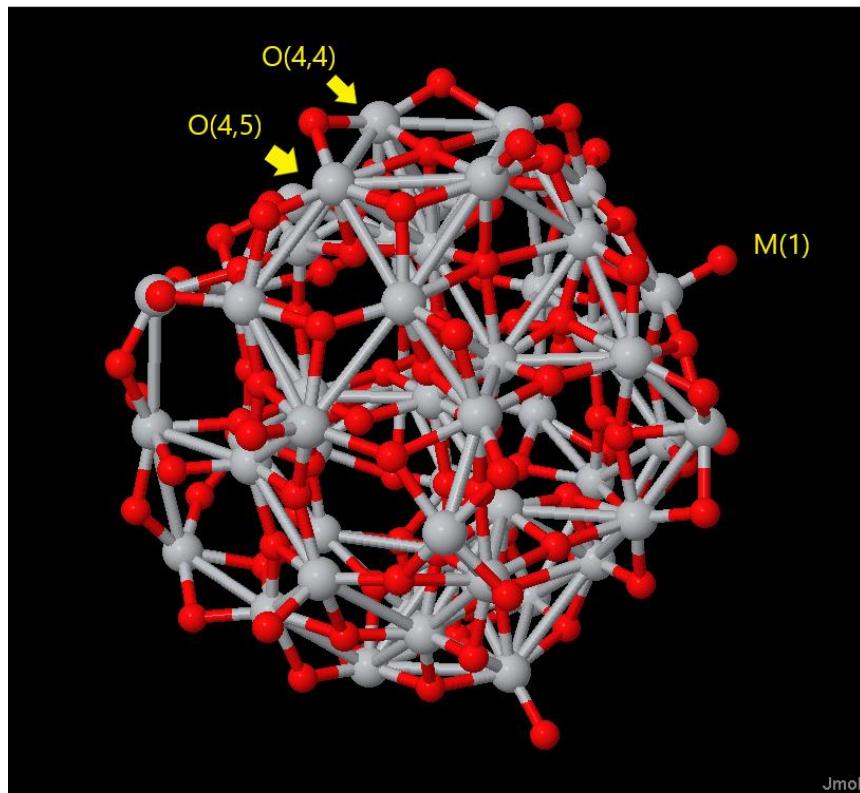


Figure 4. Nanocompound representation example

Multiple substructures have been highlighted using molviewer<sup>3</sup> in Matlab. It is seen two different configurations where an oxygen atom is a central node connecting with both metal atoms and other oxygen atoms. At the same time, we can observe multiple M (1), indicating a metal atom with a single bond.

### 3. Methodology

#### 3.1. Atena Web Portal

A publicly accessible web portal was previously developed by a Software Engineer, allowing users to generate a shortened version of the NanoFingerprint, which contained only Section 1 and 2 (described above). Given the website was still in preliminary stages of development, a preliminary assessment was needed to understand the requirements and modifications to be included, more considering the feasibility of adding them given the constraints associated

---

<sup>3</sup> <https://uk.mathworks.com/help/bioinfo/ref/molviewer.html>

with the system hosting the website itself. Site is currently accessible through [atena.urv.cat/model/](http://atena.urv.cat/model/).

Focus was put on developing the backend algorithm running the website. Since initial system was written in Go language, it was considered whether using a different programming language would have an impact on performance. As seen in [8], Golang does make better use of CPU and memory available, which is linked to the fact of compiling it to a single binary file, compared to NodeJs. Additionally, including modifications on top of the initial system in Go reduced the time needed to include all requirements.

### 3.1.1. Characteristics

As stated above, web portal backend has been developed Golang App, together with Fiber framework<sup>4</sup>. Frontend consists of Bootstrap/JS as a Single Page Application. Compiled file for website is currently hosted in a Linux Virtual Machine in Atena host. Modifications applied to code have to be performed from a separate device labelled as Aura. This configuration was established to restrict outbound traffic from Atena, limiting only SSH as well as HTTPS inbound traffic. Guidelines on how to operate both Aura and Atena are included in Appendix A. Structure of the project goes as seen below:

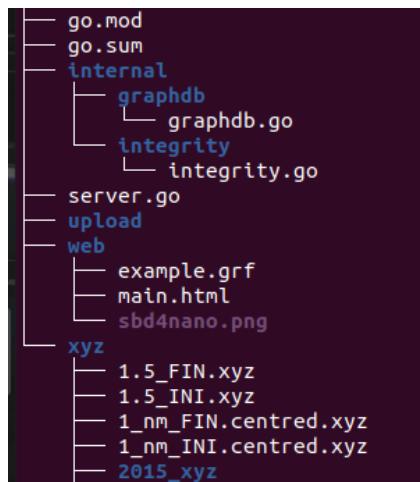


Figure 5. Project structure

*Graphdb.go* is in charge of performing all the required calculations to compute the toxicity level given a set of parameters, as well as generating the NanoFingerprint from a given nanocompound, using Neo4j<sup>5</sup> to perform a graph analysis of the nanocompound.

---

<sup>4</sup> <https://gofiber.io/>

<sup>5</sup> <https://neo4j.com/>

*Server.go* works as the main file of the application, in charge of starting the server and handling parameters submitted by user as well as providing the results through the web portal.

*Main.html* corresponds to the frontend part of the web portal. As stated earlier, it is an SPA. In terms of maintenance and update of the web portal, this is the preferable approach, more considering the different functionalities available.

### 3.1.2. Development environment

Given site was already accessible by external users, no modification was performed directly in the server hosting the site. Instead, a local testing environment was used to visualise the impact of the additional features as well as modifications performed.

A local Linux virtual machine was built from the scratch on a Windows 10 host using Oracle VM VirtualBox<sup>6</sup>, replicating the structure already in place in production. Initially the idea was to include the hardware constraints present in production, however this led to higher processing time when handling some of the nanocompounds for further analysis. Below are hardware details from the host device in the local environment:

Hostname	DESKTOP-4K5CDJT
Processor	AMD Ryzen7-3800X 8-Core Processor 3.90 GHz
Graphic Card	AMD RX 6700 XT
RAM Available	24.0 GB

Out of the available resources, around 18GB RAM have been allocated for the Linux VM, as well as using five out of all the processors available. These values have been fine-tuned as new NanoFingerprints were generated, given some of these could not be completed due to not having enough memory available. Below a comparison between testing environment and production environment:

	PROD Environment	DEV Environment
# Processors	1	5
RAM	4 GB	18GB

---

<sup>6</sup> <https://www.virtualbox.org/wiki/VirtualBox>

### 3.2. Webpage sections

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction NanoFingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

You have reached the front page of NanoFingerprint. We provide three different tools applied to nanocompounds: NanoFingerprint generation, Toxicity prediction and subcomponent searching.

**NanoFingerprint generation:**  
Given a nanocompound, a shell thickness and the maximum number of bonds per atom that are going to be represented in the NanoFingerprint, the system returns the NanoFingerprint of its shell and also the composition of the shell. The NanoFingerprint is a signature of the nanocompound based on frequency distribution scale of some specific small combinations of connected atoms.

**Toxicity prediction:**  
We provide some equations to predict toxicity of nanocompounds extracted from several references.  
You must introduce the required parameters for each equation and the system returns the toxicity level.

**Subcomponent searching:**  
Given a nanocompound, a shell thickness and a subcomponent, the system returns the atoms in the nanocompound shell that are mapped to the subcomponent.

[Nanocompound documentation](#)

Figure 6. Atena web portal homepage

When accessing web portal, users are presented with a set of tabs that allows to gather different metrics as well as generation of the NanoFingerprint. Sections below include a description of the different tabs present as well as the functionalities included from the final user point of view.

#### 3.2.1. NanoFingerprint

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction NanoFingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

**NanoFingerprint Generation**

You must provide the nanocompound (in ".xyz" format), the thickness of the shell (in nanometres) and the maximum number of bounds considered in the NanoFingerprint. Our system computes four files: the NanoFingerprint of the Shell in a ".txt" file in two formats (verbose and no verbose version)\* and also the composition of the shell (in a ".xyz" format) in two different files: one with only the atoms of the shell and the second one with the whole nanocompound in which the atoms of the shell have been labelled with "#Shell".  
Reference: Francesc Serratosa, Susana Alvarez, Laura Escorihuela and Monica Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoWeek & NanoCommons Final Conference, Cyprus 2022.

Thickness of the shell (nm)

Maximum number of bounds per atom

Structure file (in XYZ format and max size:500 Kb)

Browse... No file selected.

Calculate Reset

Figure 7. NanoFingerprint generation main page

This section allows the generation of the NanoFingerprint given a nanocompound. User is required to load a .xyz file and provide a thickness as well as a maximum number of bonds that are desired to be analysed. If user fails to input the data in the required format, portal will notify of the error and request the user to provide corrected information. There is a limitation in file size, directly tied to the capabilities of the system. These constraints are not present in the testing environment, as the available resources provided better manoeuvrability handling nanocompounds with higher file size. As such, production web portal is limited to a file size of 500Kb due to hardware constraints, whereas testing environment allows to increase this limit to 10Mb.

The screenshot shows a user interface for generating a NanoFingerprint. It includes fields for 'Thickness of the shell (nm)' (set to 0) and 'Maximum number of bounds per atom' (set to 0), both of which have red validation error messages: 'Value must be between 0.1 and 1000' and 'Value must be between 4 and 15'. Below these is a field for 'Structure file (in XYZ format and max size:500 Kb)' with a 'Browse...' button and a placeholder 'No file selected.' A red error message 'Must provide a .xyz text file.' is displayed below the input field. At the bottom are 'Calculate' and 'Reset' buttons.

*Figure 8. Error catch example*

The screenshot shows a user interface for generating a NanoFingerprint with successful input. The 'Thickness of the shell (nm)' is set to 0.4 and 'Maximum number of bounds per atom' is set to 8. The 'Structure file (in XYZ format and max size:500 Kb)' field contains 'tio2\_030.xyz' and is highlighted in grey. At the bottom are 'Calculate' and 'Reset' buttons, and three download links: 'Download Verbose NanoFingerprint', 'Download NanoFingerprint', and 'Download XYZ Compressed (Shell & Comments)'.

*Figure 9. Results availability example*

From a backend perspective, the file is locally stored in the server and multiple transactions are performed through Neo4j, not only by adding the atoms as nodes to the graph dB, but also when finding the relationship between nodes and finding the different subgraph based on a set of relationship between atoms, as seen in [1]. These are locally stored using UUID package, which avoids conflict when generating new files as uuid is used to name the directory where these files will be located.

Once all the required operations are performed, user is prompted with 3 downloadable files: Verbose NanoFingerprint, NanoFingerprint, and XYZ Compressed. First file includes a list of non-zero values associated with the NanoFingerprint, granting an overview of the subgraph relationship found given the input file. The second file provides the whole NanoFingerprint, including all zero values as well as the values listed in the Verbose file.

Both files are provided as txt files. Finally, the third available file is a .zip file containing both the input nanocompound .XYZ file (with only those atoms that are part of the shell), and the same input file labelling those atoms that were identified as part of the shell.

```

1. Shell thickness in Angstrom:4.00000
2. Maximum number of bounds per atom:8
3. Size in Angstrom:30.069616
4. Atomic number of the metal:22
5. Oxygen atoms:505
6. Metal atoms:228
7. O(1):58
8. O(2):165
9. O(3):282
10. Ti(2):4
11. Ti(3):26
12. Ti(4):38
13. Ti(5):42
14. Ti(6):118
15. O(0,1):58
16. O(0,2):165
17. O(0,3):282
18. Ti(2,0):4
19. Ti(3,0):26
20. Ti(4,0):38
21. Ti(5,0):42
22. Ti(6,0):118
23. 13426. O(0,1)-Ti(4,0):2
24. 13435. O(0,1)-Ti(5,0):28
25. 13444. O(0,1)-Ti(6,0):28
26. 13489. O(0,2)-Ti(2,0):4
27. 13498. O(0,2)-Ti(3,0):32
28. 13507. O(0,2)-Ti(4,0):64
29. 13516. O(0,2)-Ti(5,0):80
30. 13525. O(0,2)-Ti(6,0):150
31. 13570. O(0,3)-Ti(2,0):4
32. 13579. O(0,3)-Ti(3,0):46
33. 13588. O(0,3)-Ti(4,0):86
34. 13597. O(0,3)-Ti(5,0):102
35. 13606. O(0,3)-Ti(6,0):438

```

Figure 10. Verbose NanoFingerprint example

```

1. 4.00000
2. 8
3. 30.069616
4. 22
5. 505
6. 228
7. 0
8. 58
9. 165
10. 282
11. 0
12. 0
13. 0
14. 0
15. 0
16. 0
17. 0
18. 4
19. 26
20. 38
21. 42
22. 118
23. 0
24. 0
25. 0
26. 58
27. 165
28. 282
29. 0
30. 0
31. 0

```

Figure 11. Nanofingerprint example

```

1 733
2 TV:
3 O -3.790000 -13.267900 0.398662
4 O 5.682100 9.472100 -2.772712
5 Ti 7.580000 11.370000 0.000000
6 Ti 7.580000 -3.790000 -9.510000
7 O 9.472100 3.790000 5.155962
8 Ti -1.897900 -7.580000 -11.884050
9 O -3.790000 11.370000 7.534612
10 Ti 3.790000 -9.477900 -7.131350
11 Ti 0.000000 -11.370000 0.000000
12 O -9.477900 -7.580000 5.155962
13 O -9.477900 -3.790000 9.115938
14 Ti -3.790000 -9.477900 -7.131350
15 Ti 9.472100 5.682100 -4.752700
16 Ti 1.892100 -11.370000 -2.374050
17 O 11.370000 -7.580000 -1.975388
18 O -7.580000 5.682100 -9.111338
19 Ti 9.472100 -7.580000 7.135950
20 O 7.580000 11.370000 1.979988
21 Ti 1.892100 -11.370000 7.135950
22 Ti 7.580000 7.580000 -9.510000
23 O -11.370000 -5.687900 4.358638
24 Ti -3.790000 -1.897900 11.888650
25 O 3.790000 -5.687900 9.908662
26 O 9.472100 5.682100 -2.772712
27 O 9.472100 -1.897900 -6.732688
28 O -9.477900 -1.897900 -6.732688
29 Ti -11.370000 7.580000 0.000000
30 O -7.580000 5.682100 9.908662
31 O -7.580000 3.790000 -7.530012
32 O -1.897900 -7.580000 9.115938
33 O 3.790000 -7.580000 -11.485388
34 O -9.477900 -7.580000 -0.394062
35 O -11.370000 7.580000 -1.975388
36 O -3.790000 -5.687900 9.908662
37 O 1.892100 7.580000 9.115938

```

Line 11, Column 38      Tab Size: 4      Plain Text

Figure 12. XYZ file with shell-only atoms

```

1 1776
2 TV:
3 O -3.790000 -13.267900 0.398662 #Shell
4 O 5.682100 9.472100 -2.772712 #Shell
5 Ti 7.580000 11.370000 0.000000 #Shell
6 O 0.000000 7.580000 7.534612
7 Ti 7.580000 -3.790000 -9.510000 #Shell
8 O 9.472100 3.790000 5.155962 #Shell
9 Ti -1.897900 -7.580000 -11.884050 #Shell
10 O -3.790000 11.370000 7.534612 #Shell
11 Ti 3.790000 -9.477900 -7.131350 #Shell
12 Ti 3.790000 7.580000 0.000000 #Shell
13 Ti 0.000000 -11.370000 0.000000 #Shell
14 O 1.892100 -3.790000 9.115938
15 O -9.477900 -7.580000 5.155962 #Shell
16 O -9.477900 -3.790000 9.115938 #Shell
17 Ti -3.790000 -9.477900 -7.131350 #Shell
18 Ti -3.790000 9.472100 2.378650
19 Ti 9.472100 5.682100 -4.752700 #Shell
20 O -7.580000 -3.790000 -1.975388
21 Ti 1.892100 -11.370000 -2.374050 #Shell
22 Ti 0.000000 5.682100 -7.131350
23 Ti 9.472100 -3.790000 -2.374050
24 O 11.370000 -7.580000 -1.975388 #Shell
25 O -7.580000 5.682100 -9.111338 #Shell
26 Ti -7.580000 5.682100 2.378650
27 Ti 9.472100 -7.580000 7.135950 #Shell
28 O 7.580000 11.370000 1.979988 #Shell
29 O -7.580000 0.000000 -7.530012
30 O -5.687900 1.892100 -6.732688
31 Ti 1.892100 -11.370000 7.135950 #Shell
32 Ti 7.580000 7.580000 -9.510000 #Shell
33 O -11.370000 -5.687900 4.358638 #Shell
34 Ti -3.790000 -1.897900 11.888650 #Shell
35 O 3.790000 -5.687900 4.358638
36 O 3.790000 -5.687900 9.908662 #Shell
37 O -7.580000 5.682100 0.398662

```

Line 12, Column 36      Tab Size: 4      Plain Text

Figure 13. XYZ file with labels indicating shell atoms.

### 3.2.2. Toxicity Prediction

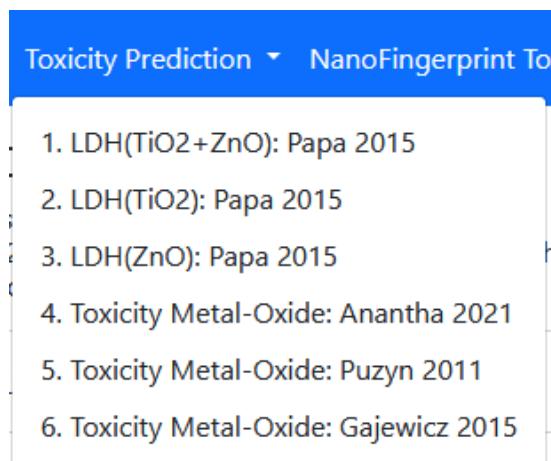


Figure 14. Toxicity Prediction

Within the tab there is a total of six separate options to be considered. The parameters required change depending on the model evaluated. As seen in [3], three different equations are presented, which have been added as individual fields since the input parameters differ between them.

### Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction → NanoFingerprint Toxicity Prediction → Subcomponent Search Examples Data Input Format About us

**LDH(TiO<sub>2</sub>+ZnO)**

This equation predicts the Lactate Dehydrogenase Release (LDH) of TiO<sub>2</sub> and ZnO. It corresponds to Equation 1 in the following reference:E. Papa, J.P. Doucet & A. Doucet-Panaye (2015): Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors, SAR and QSAR in Environmental Research, DOI: 10.1080/1062936X.2015.1080186

Concentration (mg/L)	Size (diameter in nm)
Size in Pbs (diameter in nm)	
<input type="button" value="Calculate"/>	<input type="button" value="Reset"/>

Figure 15. LDH (TiO<sub>2</sub> + ZnO) toxicity calculation

### Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction → NanoFingerprint Toxicity Prediction → Subcomponent Search Examples Data Input Format About us

**LDH(TiO<sub>2</sub>)**

This equation predicts the Lactate Dehydrogenase Release (LDH) of TiO<sub>2</sub>. It corresponds to Equation 2 in the following reference:E. Papa, J.P. Doucet & A. Doucet-Panaye (2015): Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors, SAR and QSAR in Environmental Research, DOI: 10.1080/1062936X.2015.1080186

Concentration (mg/L)	Size (diameter in nm)
<input type="button" value="Calculate"/>	<input type="button" value="Reset"/>

Figure 16. LDH (TiO<sub>2</sub>) toxicity calculation

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction Nanofingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

### LDH(ZnO)

This equation predicts the Lactate Dehydrogenase Release (LDH) of ZnO. It corresponds to Equation 3 in the following reference:E. Papa, J.P. Doucet & A. Doucet-Panaye (2015): Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors, SAR and QSAR in Environmental Research, DOI:10.1080/1062936X.2015.1080186

Concentration (mg/L)	Size in Water (diameter in nm)
Size in Pbs (diameter in nm)	

**Calculate** **Reset**

Figure 17. LDH (ZnO) toxicity calculation

These values are limited to fall within a given interval, if there is a failure to input the correct value, no toxicity will be calculated, and user will be prompted to correct the input values, as seen in the example below.

### LDH(TiO<sub>2</sub>+ZnO)

This equation predicts the Lactate Dehydrogenase Release (LDH) of TiO<sub>2</sub> and ZnO. It corresponds to Equation 1 in the following reference:E. Papa, J.P. Doucet & A. Doucet-Panaye (2015): Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors, SAR and QSAR in Environmental Research, DOI: 10.1080/1062936X.2015.1080186

Concentration (mg/L) 0	Size (diameter in nm) 0
Value must be between 6.25 and 800	
Size in Pbs (diameter in nm) 0	
Value must be between 39.5 and 15484	

**Calculate** **Reset**

Figure 18. Example of incorrect input parameters submission

### LDH(TiO<sub>2</sub>)

This equation predicts the Lactate Dehydrogenase Release (LDH) of TiO<sub>2</sub>. It corresponds to Equation 2 in the following reference:E. Papa, J.P. Doucet & A. Doucet-Panaye (2015): Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors, SAR and QSAR in Environmental Research, DOI: 10.1080/1062936X.2015.1080186

Concentration (mg/L) 25	Size (diameter in nm) 30
<b>Calculate</b> <b>Reset</b>	
Ldh 0.79	

Figure 19. Example of output for LDH (TiO<sub>2</sub>)

# Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction NanoFingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

## Toxicity Metal-Oxide: Anantha 2021

Prediction of toxicity of Metal-Oxide (binary output). It corresponds to the model presented in:  
NanoTox: Development of a Parsimonious In Silico Model for Toxicity Assessment of Metal-Oxide Nanoparticles Using Physicochemical Features, Nilesh Anantha Subramanian and Ashok Palaniappan, ACS Omega 2021 6 (17), 11729-11739, DOI: 10.1021/acsomega.1c01076.

Size (nm)	HydroSize (nm)
SurfCharge (eV)	SurfArea (nm <sup>2</sup> )
E <sub>c</sub> (eV)	Time (s)
Dose (mg)	Eneg (eV)
NOxygen	

**Calculate** **Reset**

Figure 20. Anantha toxicity prediction

Anantha section does not provide a numeric value as it would occur with the other models, but rather a statement indicating whether the analysed nanocompound is either toxic or not, as seen in the example below:

## Toxicity Metal-Oxide: Anantha 2021

Prediction of toxicity of Metal-Oxide (binary output). It corresponds to the model presented in:  
NanoTox: Development of a Parsimonious In Silico Model for Toxicity Assessment of Metal-Oxide Nanoparticles Using Physicochemical Features, Nilesh Anantha Subramanian and Ashok Palaniappan, ACS Omega 2021 6 (17), 11729-11739, DOI: 10.1021/acsomega.1c01076.

Size (nm) 39.7	HydroSize (nm) 267
SurfCharge (eV) 36.3	SurfArea (nm <sup>2</sup> ) 64.7
E <sub>c</sub> (eV) -1.51	Time (s) 24
Dose (mg) 0.001	Eneg (eV) 1.61
NOxygen 3	

**Calculate** **Reset**

It's NO toxic



Figure 21. Anantha prediction example

Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction NanoFingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

Toxicity Metal-Oxide: Puzyn 2011

Prediction of toxicity of Metal-Oxide ( $\log(1/\text{EC50})$ ). It corresponds to the model presented in:  
Using nano-QSAR to predict the cytotoxicity of metal oxide nanoparticles, Tomasz Puzyn, Bakhtiyor Rasulev, Agnieszka Gajewicz, Xiaoke Hu, Thabitah P Dasari, Andrea Michalkova, Huey-Min Hwang, Andrey Toropov, Danuta Leszczynska, Jerzy Leszczynski, Nature Nanotechnology, 2011 M6(3), DOI: 10.1038/nano.2011.10.

Calculate
Reset

Figure 22. Toxicity prediction using equation from Puzyn

Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction NanoFingerprint Toxicity Prediction Subcomponent Search Examples Data Input Format About us

Toxicity Metal-Oxide: Gajewicz 2015

Prediction of toxicity of Metal-Oxide ( $\log(1/\text{EC50})$ ). It corresponds to the model presented in:  
Agnieszka Gajewicz, Nicole Schaeublin, Bakhtiyor Rasulev, Saber Hussain, Danuta Leszczynska, Tomasz Puzyn & Jerzy Leszczynski (2015) Towards understanding mechanisms governing cytotoxicity of metal oxides nanoparticles: Hints from nano-QSAR studies, Nanotoxicology, 9:3, 313-325, DOI: 10.3109/17435390.2014.930195

Calculate
Reset

Figure 23. Toxicity prediction using equation from Gajewicz

In the backend, different REST API calls are performed depending on the selected options. Each call will perform a different operation considering input values. These operations are a code representation of the equations presented in [4] [3] [5] [6]. Example of the code used for computing the toxicity levels is appended below. Once computed these are returned to the webportal in JSON format, which then is:

```

147 func ldh4tz(c *fiber.Ctx) error {
148     input := l4tz{}
149     if e := c.BodyParser(&input); e != nil {
150         return status400(c, e)
151     }
152     if e := integrity.CheckFields(input, "l4tz"); e != "" {
153         return c.JSON(ldh{Ldh: 0, Error: e})
154     }
155     log.Println(c.Path(), input.Concentration, input.Size, input.SizeInPbs)
156
157     return c.JSON(ldh{Ldh: roundFloat(0.66+0.003*input.Concentration+0.005*input.Size-5*input.SizeInPbs, 2), Error: ""})
158 }
159 }
```

Figure 24. Toxicity prediction for ( $TiO_2+ZnO$ ) model

```

160 func ldh4tio2(c *fiber.Ctx) error {
161     input := l4t{}
162     if e := c.BodyParser(&input); e != nil {
163         return status400(c, e)
164     }
165     if e := integrity.CheckFields(input, "l4t"); e != "" {
166         return c.JSON(ldh{Ldh: 0, Error: e})
167     }
168     log.Println(c.Path(), input.Concentration, input.Size)
169
170     return c.JSON(ldh{Ldh: roundFloat(0.599+0.003*input.Concentration+0.004*input.Size, 2), Error: ""})
171 }
172
173 func ldh4zno(c *fiber.Ctx) error {
174     input := l4z{}
175     if e := c.BodyParser(&input); e != nil {
176         return status400(c, e)
177     }
178     if e := integrity.CheckFields(input, "l4z"); e != "" {
179         return c.JSON(ldh{Ldh: 0, Error: e})
180     }
181
182     log.Println(c.Path(), input.Concentration, input.Size, input.Pbs)
183
184     return c.JSON(ldh{Ldh: roundFloat(1.041+0.001*input.Size-0.001*input.Pbs+0.001*input.Concentration, 2), Error: ""})
185 }
186

```

Figure 25. Toxicity models for TiO<sub>2</sub> and ZnO

### 3.2.3. NanoFingerprint Toxicity Prediction

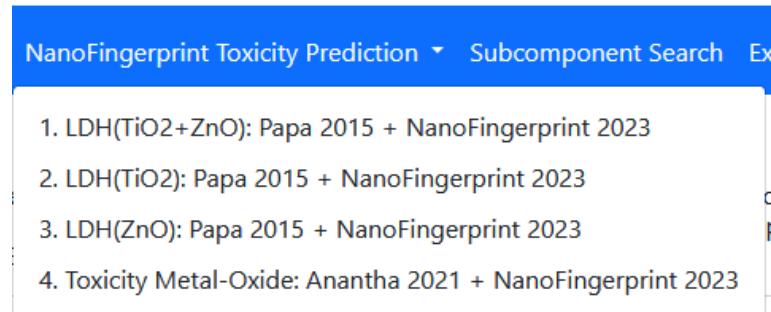


Figure 26. NanoFingerprint toxicity prediction options

In this scenario, toxicity levels are computed using not only the parameters previously presented but adding as well a previous computed NanoFingerprint. 4 different options are presented below, being 3 of these linked to [3] equations.

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction ▾ NanoFingerprint Toxicity Prediction ▾ Subcomponent Search Examples Data Input Format About us

**1. LDH(TiO<sub>2</sub>+ZnO): Papa 2015 + NanoFingerprint 2023**

Prediction of toxicity of Metal-Oxide (LDH). It corresponds to the model presented in:  
F. Serratosa, S. Álvarez, L. Escorihuela and M. Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoWeek & NanoCommons Final Conference 2022, Cyprus 2022.

NanoFingerprint (Thickness: 0.4nm, MaxBonds: 8)

No file selected.

Figure 27. TiO<sub>2</sub> + ZnO + NanoFingerprint

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction ▾ NanoFingerprint Toxicity Prediction ▾ Subcomponent Search Examples Data Input Format About us

**LDH(TiO<sub>2</sub>): Papa 2015 + NanoFingerprint 2023**

Prediction of toxicity of Metal-Oxide (LDH). It corresponds to the model presented in:  
F. Serratosa, S. Álvarez, L. Escorihuela and M. Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoWeek & NanoCommons Final Conference 2022, Cyprus 2022.

No file selected.

Figure 28. TiO<sub>2</sub> + NanoFingerprint

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction ▾ NanoFingerprint Toxicity Prediction ▾ Subcomponent Search Examples Data Input Format About us

**1. LDH(ZnO): Papa 2015 + NanoFingerprint 2023**

Prediction of toxicity of Metal-Oxide (LDH). It corresponds to the model presented in:  
F. Serratosa, S. Álvarez, L. Escorihuela and M. Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoWeek & NanoCommons Final Conference 2022, Cyprus 2022.

NanoFingerprint (Thickness: 0.4nm, MaxBonds: 8)

No file selected.

Figure 29. ZnO + NanoFingerprint

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction ▾ NanoFingerprint Toxicity Prediction ▾ Subcomponent Search Examples Data Input Format About us

### 4. Toxicity Metal-Oxide: Anantha 2021 + NanoFingerprint 2023

Prediction of toxicity of Metal-Oxide (binary output). It corresponds to the model presented in:  
F. Serratosa, S. Álvarez, L. Escorihuela and M. Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoWeek & NanoCommons Final Conference 2022, Cyprus 2022.

Size (nm)	HydroSize (nm)
SurfCharge (eV)	SurfArea (nm <sup>2</sup> )
Ec (eV)	Time (s)
Dose (mg)	Eneg (eV)
NOxygen	
NanoFingerprint (Thickness: 0.4nm, MaxBonds: 8)	
Browse...	No file selected.

**Calculate** **Reset**

Figure 30. Anantha equation + NanoFingerprint

Pre-processing techniques are applied to the input NanoFingerprint before applying optimal models for toxicity prediction. Additional details regarding these models as well as the pre-processing techniques used can be seen in Results section. Like the previous tab description, all models except “Anantha” provide a toxicity estimation.

### 3.2.4. Subcomponent Search

## Nanocompound analyser based on sub-structure extraction

Home NanoFingerprint Toxicity Prediction ▾ NanoFingerprint Toxicity Prediction ▾ Subcomponent Search Examples Data Input Format About us

### Subcomponent Search

You must provide the nanocompound in ".xyz", the substructure in ".xyz" or ".grf" format and the shell thickness in nanometres. The system obtains the atoms in the nanocompound that are isomorphic to the subcomponent that are located at its shell, this information is saved in two different files (format '.xyz'): one with the subcomponents found in the shell and the second is the nanocompound in which the atoms in the found subcomponents in the shell have been labelled which #Subcomponent.

Reference: Francesc Serratosa, Susana Álvarez, Laura Escorihuela and Monica Calatayud, Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration. NanoCommons Final Conference, Cyprus 2022.

The subgraph matching algorithm is VF3 (Challenging the time complexity of exact subgraph isomorphism for huge and dense graphs with VF3 - Carletti V., Foglia P., Saggesse A., Vento M. - IEEE transactions on pattern analysis and machine intelligence - 2018) and we use the implementation in <https://github.com/MiMaLab/vf3lib>

Nanocompound file (in XYZ format)	Subcomponent file (in XYZ or GRF format)
Browse... No file selected.	Browse... No file selected.
Thickness of the shell (nm)	
<b>Calculate</b>	<b>Reset</b>

Figure 31. Subcomponent search section

Section allows the user to find patterns in the shell of the nanoncompound within a given nanocompound and a substructure. Function for this computation is not included in

graphdb.go but rather in a set of multiple files written in C and based on VF3<sup>7</sup>. Server.go oversees the execution of the C code through the exec library.

The screenshot shows a user interface for a nanocompound analyser. It has two main sections: 'Nanocompound file (in XYZ format)' and 'Subcomponent file (in XYZ or GRF format)'. In the first section, a file 'Al2O3\_013.xyz' is selected. In the second section, a file 'Al2O3\_01\_01.grf' is selected. Below these, a 'Thickness of the shell (nm)' input field contains '0.5'. A 'Calculate' button is visible. Underneath, a 'Number of detected atoms' input field shows '24'. At the bottom, a green 'Download the output file' button is present.

Figure 32. Substructure search result example

### 3.2.5. Examples

The screenshot shows a website for the Nanocompound analyser. The header includes links for Home, NanoFingerprint, Toxicity Prediction, NanoFingerprint Toxicity Prediction, Subcomponent Search, Examples, Data Input Format, and About us. The main content area is titled 'Nanocompound analyser based on sub-structure extraction'. Under the 'Examples' heading, there are two buttons: 'Download modelling examples' (grey) and 'Download subcomponent search examples' (green). Below these buttons are three logos: 'Flame', 'nanoplanoATX', and 'SbD Nano4'.

Users are provided with different samples to test both “NanoFingerprint” and “Subcomponent Search” tabs. Modelling examples files contain both images of some nanocompounds and XYZ files linked to the same nanocompounds to be used as input. At the same time, the corresponding NanoFingerprint has been added as part of the output examples. The second file contains both nanocompounds as well as substructures associated to these (in .GRF format). Some output examples are also added as a demonstration of the expected results.

<sup>7</sup> <https://github.com/MiviaLab/vf3lib>

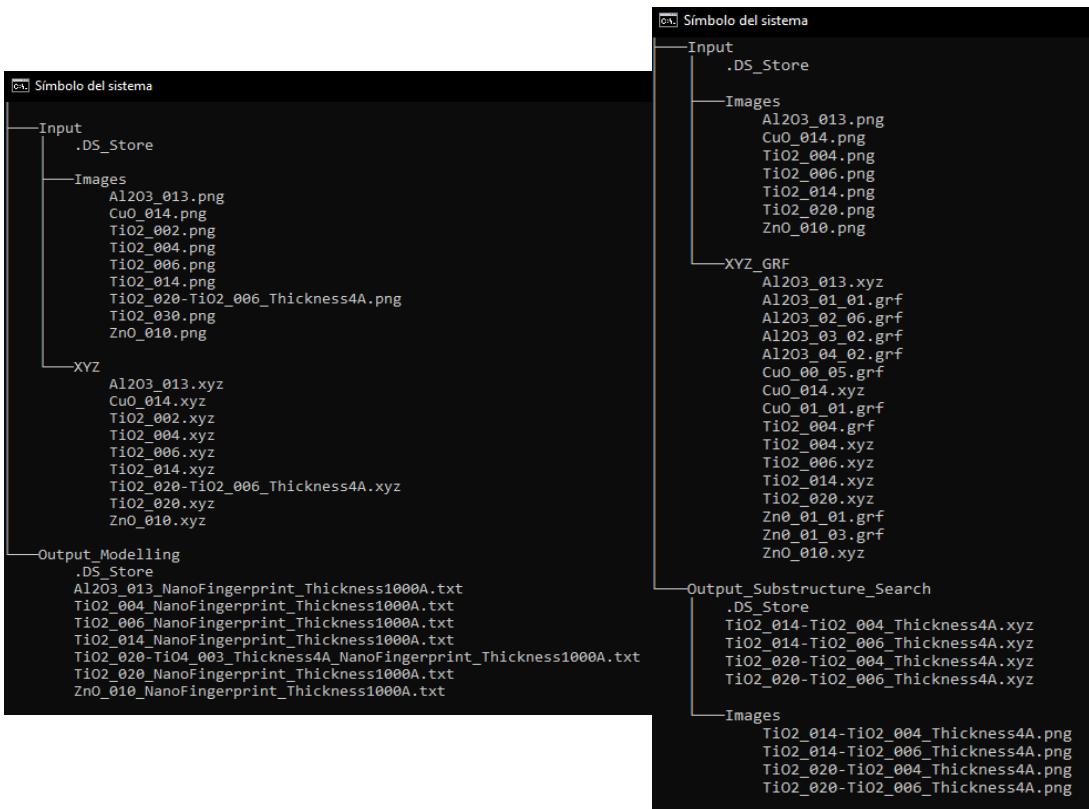


Figure 33. Directory tree for the Modelling and Subcomponent examples

### 3.2.6. Data Input Format

#### Nanocompound analyser based on sub-structure extraction

[Home](#) [NanoFingerprint](#) [Toxicity Prediction](#) ▾ [NanoFingerprint Toxicity Prediction](#) ▾ [Subcomponent Search](#) [Examples](#) [Data Input Format](#) [About us](#)

##### Data Input Format

These are the data input formats required in most of the programs provided.

.XYZ:

[https://www.wikiwand.com/en/XYZ\\_file\\_format](https://www.wikiwand.com/en/XYZ_file_format)

.GRF:

On the first line there must be the number of atoms. The next lines contain the atomic number of the atoms, one atom per line, preceded by the atom id; no atom ids must be in the range from 0 to the number of atoms - 1. Then, per each atom there is the number of bonds coming out of the atom, followed by a line for each bond containing the ids of the edge ends.

[Download an example file of a Titanium connected to 6 Oxygen.](#)



Figure 34. Expected input files details.

### 3.3. Nanofingerprint Generation

As stated earlier, graphdb.go file is in charge of performing the required operations to generate the NanoFingerprint given a set of input parameters. Once the request is sent through the webportal, server.go captures the values from the form and performs a go routine call to another function (coVectorCalc), which runs in parallel to main thread.

Within this function, a call to the main function in graphdb.go file, named GenVector, is performed.

```
func vectorCalc(c *fiber.Ctx) error {
    thickness := c.FormValue("thickness")
    bounds := c.FormValue("bounds")
    file, _ := c.FormFile("file")
    strId := uuid.New().String()
    // fileName := fmt.Sprintf(uploadDir+"%s", strId+".xyz")
    filename_form := c.FormValue("filename")
    fileName := fmt.Sprintf(uploadDir+"%s", filename_form)
    c.SaveFile(file, fileName)
    go coVectorCalc(thickness, bounds, fileName, strId, filename_form)

    return c.JSON(modId{ID: strId, Error: ""})
}
```

Figure 35. Form input variables parsing

```
func coVectorCalc(thickness string, bounds string, fileName string, strId string, base_name string) {
    var thick float64
    var bo int
    var err error
    if thick, err = strconv.ParseFloat(thickness, 64); err != nil {
        log.Panic(err)
    }
    if bo, err = strconv.Atoi(bounds); err != nil {
        log.Panic(err)
    }
    fileToDownload, fileToDownloadNoVer, err, fileShellZIP := graphdb.GenVector(fileName, thick, false, bo, downloadDir, strId, base_name)
    lock.Lock()
    defer lock.Unlock()
    if err == nil {
        fpReq[strId] = mod{File: fileToDownload, FileNoVer: fileToDownloadNoVer, Error: "", FileShell: fileShellZIP}
    } else {
        fpReq[strId] = mod{File: "", FileNoVer: "", Error: "Processing of the data failed"}
    }
}
```

Figure 36. Go routine call to CoVectorCalc

GenVector reads the input file submitted through the form and parses the contents of the file to compute the distance between atoms, as well as determining whether the atom corresponds to the shell part of the nanocompound or not. From this function, a call to computeShell is performed, which creates a Neo4j session. During the session, the atoms identified as shell are added to a graph database and queries are submitted to Neo4j to find those nodes that match the specifications listed in Sections 3 and 4 from [1]. Once the required operations are performed, Neo4j session is closed and database is deleted, to ensure no remaining data is left behind that could impact future NanoFingerprint generations. Results from Neo4j are written to 2 different files, one containing the whole NanoFingerprint array, and a second file specifying only those non-zero values found during the analysis. Furthermore, 2 more files are generated, which are stored as a ZIP file containing details regarding the shell itself, in .XYZ format. Details on GenVector function as well as secondary functions within graphdb.go are included in Appendix A.

### 3.3. Dataset

For this project, 2 different datasets have been used, one using the samples presented in [3] and adding the NanoFingerprints associated with the nanocompounds listed in the dataset, and a second dataset using the parameters included in [4], adding as well the NanoFingerprints for these nanocompounds. Regarding the first dataset, which will be labelled as “Papa\_db” for simplicity when referring to it on this document, is formed by a total of forty-two samples. Out of these, nine different nanocompounds were identified.

Second dataset, which will be referred as “Anantha\_db” going forward on this document, is formed by twenty-one different nanocompounds and a total of 483 samples. Both datasets have NanoFingerprints with thickness of 0.4nm and maximum number of bonds set to 8 as default values. Considering the bigger the value of MaxBonds the lower the computational efficiency, no information loss was observed when using a value of 8, as no substructure was found with more than 8 bonds. Most of the samples used for the toxicity prediction in both datasets showed a maximum of 6 bonds with other metal or oxygen atoms. Regarding thickness, it was set to 0.4nm as it was the most common value seen on the nanocompounds analysed. Thickness reference was initially extracted from the “Nanoparticle shell-depth calculator” portal<sup>8</sup>. Given these parameters, the NanoFingerprint array length was equivalent to 19867, as seen in [1]. This implies that the size of Papa\_db is [42, 19874], whereas Anatha\_db is [483, 19893]. These sizes correspond to the raw data without any preprocessing techniques applied. Since the number of samples available in Papa\_db was low, cross-validation was performed to analyse the features as an attempt to avoid overfitting the model. Prior initial analysis consisted in dividing data in training, validation and test by assigning 70% of samples to training, and 15% for both validation and test and testing a set of models after performing a data-driven analysis. Cross-validation was also performed for Anantha\_db. Additionally, recursive feature elimination (RFE) has been applied. As presented in [9], this technique allows to prune features considering the performance on the model and a ranking criterion. The iterative process of reducing features that do not contribute to a better model performance allows to find and rank those features that do so, especially for a limited number of samples. Results associated to this initial approach are included in section 4.

---

<sup>8</sup> <https://nanogen.me/shell-depth>

## 3.4. Metrics

### 3.3.1. CCC

Presented in [10], metric is presented as an alternative to validate QSAR models, which are described as classification algorithms used mainly in chemical engineering. Formula for the metric is presented below:

$$CCC = \frac{2 * \sum_{i=1}^n (y_i - \mu_y) * (y_{predicted_i} - \mu_{y_{predicted}})}{\sum_{i=1}^n (y_i - \mu_y)^2 + \sum_{i=1}^n (y_{predicted_i} - \mu_{y_{predicted}})^2 + n * (\mu_y - \mu_{y_{predicted}})^2}$$

In this case,  $CCC \leq 0$  would indicate model does not fit data properly and fails to predict values for unseen data. If CCC is closer to 1, it can be confirmed with high confidence that model is a good candidate. Information would need to be correlated with other model parameters.

### 3.3.2. RMSE

Allows to measure the accuracy of a model by comparing the predicted data obtained from a model when analysing unseen data with the ground truth associated with these testing samples. The benefit of using this metric is linked to the impact large errors has on the metric, as the distance between prediction and actual value is later squared.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - y_{predicted_i})^2}{n}}$$

### 3.3.3. MAE

Another metric that has been considered for model evaluation has been “Mean Absolute Error”, which measures the average of the error in predicting the output, this gives a straightforward measure of the error size across the analysed samples. A larger MAE value indicates a poor performance of the model, since difference between prediction and expected value would be higher.

$$\frac{1}{n} \sum_{i=1}^n |y_i - y_{predicted_i}|$$

### 3.3.4. R<sup>2</sup>

The coefficient of determination, or r-squared, has been widely used in QSAR model analysis – it provides details regarding how well the model account for the variability of the

expected values. r-squared of 0 would indicate model does not account for the samples analysed.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - y_{predicted_i})^2}{\sum_{i=1}^n (y_i - \mu_y)^2}$$

## 4. Experimental Validation

### 4.1. Environment

In order to perform the analysis of the multiple datasets obtained, Python has been used. Multiple libraries were used to interact with the datasets, such as pandas and numpy. Models have been trained using scikit-learn library. Regarding the hardware features of the device used, the CPU was 8-core and a total of 24GB RAM available. In this case however, given the volume of data, no limitation was seen as per the resources used. Code used for preprocessing and modelling are included in Appendix A.

VSCode was the software using for working in the modelling script as well as the files required for the web portal, as stated in a previous section, web portal functionality was tested in a local Linux VM.

### 4.2. Web Portal Development

As stated in the definition of the project, the web portal was able to perform limited operations when generating the NanoFingerprint. Both sections 3 and 4 are seen to be the most exhaustive in terms of number of operations required, and these were not initially added to the production environment as it would disable the access to the portal itself. As part of the development of the website, these were later included with no direct impact on the performance of the server and are available to be downloaded by the end user. Additionally, user was given the option to download a new .XYZ with atoms linked to the nanocompound shell, as well as the same .XYZ but labelling what atoms were part of the shell. Both files are zipped to be downloaded as a single file. In order to avoid any possible disability of the website due to high processing time, files loaded to portal have been limited to 500Kb when calculating the NanoFingerprint.

Some additional functionalities were added linked to parsing the input file, which was seen to fail in some instances when these did not include a comment in the second line. This was corrected to ensure all files could be properly parsed as well as handling the error without crashing the web portal.

A better file management was added – initially, uuid generated from Google UUID was used to label the files that would later be downloaded. However, this added the difficulty of identifying what nanocompound it was referring to as well as what parameters were used for both thickness and maximum bonds. The proposed method consists of using uuid as a directory name instead, keeping name convention used in input file as well as adding what parameters were included, as seen in the example below:

```
vboxuser@ubuntu-desktop:~/atenaserver/downloads/39ff738d-3501-41b0-b84e-3f479e3691a9$ tree
.
├── SiO2_015_4nm_8_Nanofingerprint.txt
├── SiO2_015_4nm_8_Verbose_Nanofingerprint.txt
└── SiO2_015_compressed.zip
    └── SiO2_015_shell_labeled.xyz
        └── SiO2_015_shell.xyz

0 directories, 5 files
```

Proposed Linear Regression models have been included in webportal as alternative toxicity prediction algorithms using the information provided by the NanoFingerprints. Website has been entirely rebuilt to allocate all changes presented above.

#### 4.3. Data Preprocessing

Given the vast number of features available after appending the NanoFingerprints to the different databases, feature selection was required, specially to remove those features that had redundant zeros across all samples. Some nanocompounds were initially removed due to not having the associated XYZ file available. For Papa\_db, this was linked to the compound file size, which was around GBs of data. For Anantha\_db, some nanocompounds could not generate a NanoFingerprint due to memory constraints in both testing and production environments.

Three different subsets have been created for comparing results with the ones provided in [3], whereas Anantha\_db did not require any subset. First subset consists only of nanocompounds labelled as “TiO<sub>2</sub>”. Second subset is formed by “ZnO” components, and third one is composed by both nanocompound types. Data was initially split in training and test. When applying data-driven analysis to “TiO<sub>2</sub>” dataset, no impact in the final number of rows were seen, however size of available features did change in both “ZnO” and “TiO<sub>2</sub>+ZnO” datasets. In all instances, cross-validation together with recursive feature elimination has been performed to see as well what features would work best given the training data provided. Following sections provide an example of possible training partitions after only applying data-driven analysis and removing features that are known to impact the model negatively. These are part of a data-driven analysis and not definitive but provides a guidance on how preprocessing affect to the final remaining features. In the following

section, when presenting the models obtained, a mention on the features used as well as the results obtained for the model are included.

#### 4.3.1. Papa – TiO<sub>2</sub> dataset

Below is an example of some rows under TiO<sub>2</sub> subset:

	Name	Eng. Size (X0)	Size in Water (X1)	Size in PBS (X2)	Concentration (X4)	Zeta Potential (X5)	Y - LDH	Thickness (nm)	MaxBounds	Size (Angstrom)	Atomic Number Metal	Number Oxygen Atoms	Number Metal Atoms
2	TiO <sub>2</sub> _30	30	125	1250	25	-10	0.9	4	8	30.069616	22	505	228
3	TiO <sub>2</sub> _30	30	102	987	25	-12	1	4	8	30.069616	22	505	228
4	TiO <sub>2</sub> _30	30	281	1543	50	-15	0.75	4	8	30.069616	22	505	228
5	TiO <sub>2</sub> _30	30	101	1045	50	-9	0.7	4	8	30.069616	22	505	228
6	TiO <sub>2</sub> _30	30	299	1754	100	-11	1.04	4	8	30.069616	22	505	228
7	TiO <sub>2</sub> _30	30	134	961	100	-11	1.09	4	8	30.069616	22	505	228
8	TiO <sub>2</sub> _30	30	600	1876	200	-12	1.15	4	8	30.069616	22	505	228
9	TiO <sub>2</sub> _30	30	298	1165	200	-12	1.2	4	8	30.069616	22	505	228
10	TiO <sub>2</sub> _45	45	129	2567	25	-9	0.9	4	8	45.902255	22	1320	674
11	TiO <sub>2</sub> _45	45	129	2309	25	-10	0.85	4	8	45.902255	22	1320	674
12	TiO <sub>2</sub> _45	45	201	2431	50	-9	0.75	4	8	45.902255	22	1320	674
13	TiO <sub>2</sub> _45	45	201	2987	50	-11	0.78	4	8	45.902255	22	1320	674
14	TiO <sub>2</sub> _45	45	451	2941	100	-11	1.4	4	8	45.902255	22	1320	674

Figure 37. TiO<sub>2</sub> dataset example (Part 1)

	Number O(1)	Number O(2)	Number O(3)	Number O(4)	Number O(5)	Number O(6)	Number O(7)	Number O(8)	Number M(1)	Number M(2)	Number M(3)	Number M(4)	Number M(5)	Number M(6)	Number M(7)	Number M(8)
2	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
3	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
4	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
5	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
6	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
7	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
8	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
9	58	165	282	0	0	0	0	0	0	4	26	38	42	118	0	0
10	60	368	892	0	0	0	0	0	0	0	96	108	128	342	0	0
11	60	368	892	0	0	0	0	0	0	0	96	108	128	342	0	0
12	60	368	892	0	0	0	0	0	0	0	96	108	128	342	0	0
13	60	368	892	0	0	0	0	0	0	0	96	108	128	342	0	0
14	60	368	892	0	0	0	0	0	0	0	96	108	128	342	0	0

Figure 38. TiO<sub>2</sub> dataset example (Part 2)

When working with “TiO<sub>2</sub>”, different features were initially discarded as part of feature analysis. Initial step removed the first 4 values from Section 1 (*Shell Thickness, Maximum number of bonds, Size in Angstrom, and Atomic number of the metal*). Regarding “Size in Angstrom”, this information is already provided by X0 feature, as seen above. The other listed features had a standard deviation of 0, which would not help in the classification of training samples.

2 main approaches were taken into consideration: first approach consisted in joining all values under the same section and type, whereas second would consist of just removing those features with 0 values in all samples. Even though additional

As part of the first approach we take for instance the following example: in Section 3 one of the substructures is defined as O (1,1), which would refer to the subgraph where there is an oxygen atom as a central atom connected to another oxygen atom as well as a metal atom. The idea was to group all found subgraphs that had any similar central atom, either metal

or oxygen, by joining the results obtained for any combination of metal and oxygen atoms. This reduced the number of features linked to the NanoFingerprint to only 12 different features. Below included a snippet of the remaining DB:

	0	1	2	3	4	5	6	Section2_0	Section2_M	Section3_0	Section3_M	Section4_0	Section4_M	Section4_OM
0	30.0	125.0	1250.0	25.0	-10.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
1	30.0	102.0	987.0	25.0	-12.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
2	30.0	281.0	1543.0	50.0	-15.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
3	30.0	101.0	1045.0	50.0	-9.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
4	30.0	299.0	1754.0	100.0	-11.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
5	30.0	134.0	961.0	100.0	-11.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
6	30.0	600.0	1876.0	200.0	-12.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
7	30.0	298.0	1165.0	200.0	-12.0	505.0	228.0	505.0	228.0	505.0	228.0	0.0	0.0	1064.0
8	45.0	129.0	2567.0	25.0	-9.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
9	45.0	129.0	2309.0	25.0	-10.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
10	45.0	201.0	2431.0	50.0	-9.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
11	45.0	201.0	2987.0	50.0	-11.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
12	45.0	451.0	2941.0	100.0	-11.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
13	45.0	451.0	1934.0	100.0	-9.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
14	45.0	876.0	1965.0	200.0	-11.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
15	45.0	876.0	2109.0	200.0	-10.0	1320.0	674.0	1320.0	674.0	1320.0	674.0	0.0	0.0	3048.0
16	125.0	136.0	3215.0	25.0	-11.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
17	125.0	136.0	2667.0	25.0	-10.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
18	125.0	149.0	3782.0	50.0	-10.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
19	125.0	149.0	2144.0	50.0	-15.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
20	125.0	343.0	3871.0	100.0	-12.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
21	125.0	343.0	2890.0	100.0	-9.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
22	125.0	967.0	3813.0	200.0	-9.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0
23	125.0	967.0	2671.0	200.0	-8.0	10766.0	5400.0	10766.0	5400.0	10766.0	5400.0	0.0	0.0	24552.0

Figure 39. TiO<sub>2</sub> joined sections dataset

Notice that Section 2 and Section 3 results are identical to the total number of oxygen atoms and metal atoms. Checking the generated verbose file for TiO<sub>2</sub>\_30 for example, it is seen that the substructures found on these sections are mostly referred to either oxygen atoms bonded to only other metal atoms, or metal atoms only bonded to only oxygen atoms.

```

1 1. Shell thickness in Angstrom:4.000000
2 2. Maximum number of bounds per atom:8
3 3. Size in Angstrom:30.069616
4 4. Atomic number of the metal:22
5 5. Oxygen atoms:505
6 6. Metal atoms:228
7 7. O(1):58
8 8. O(2):165
9 9. O(3):282
10 10. Ti(2):4
11 11. Ti(3):26
12 12. Ti(4):38
13 13. Ti(5):42
14 14. Ti(6):118
15 15. O(0,1):58
16 16. O(0,2):165
17 17. O(0,3):282
18 18. Ti(2,0):4
19 19. Ti(3,0):26
20 20. Ti(4,0):38
21 21. Ti(5,0):42
22 22. Ti(6,0):118
23 23. O(0,1)-Ti(4,0):2
24 24. O(0,1)-Ti(5,0):28
25 25. O(0,1)-Ti(6,0):28
26 26. O(0,2)-Ti(2,0):4
27 27. O(0,2)-Ti(3,0):32
28 28. O(0,2)-Ti(4,0):64
29 29. O(0,2)-Ti(5,0):80
30 30. O(0,2)-Ti(6,0):150
31 31. O(0,3)-Ti(2,0):4
32 32. O(0,3)-Ti(3,0):46
33 33. O(0,3)-Ti(4,0):86
34 34. O(0,3)-Ti(5,0):102
35 35. O(0,3)-Ti(6,0):438

```

Figure 40. TiO<sub>2</sub>\_30 Verbose NanoFingerprint

From the figure above, Section 2 is represented by rows 7 – 14, whereas Section 3 is represented by rows 15 – 22. When comparing the results in both sections, it is seen that same values are obtained. Going back to the partially preprocessed dataset, both “section4\_O” and “section4\_M” are filled with zeros. Correlating this information with the example above, it is seen that substructures found for TiO<sub>2</sub> nanocompounds are mostly between oxygen and metal atoms, no substructure was found with labels O(x,y) – O(x',y') or Ti(x,y) – Ti(x',y') to be part of the shell of the nanocompound.

Removing column features for both Section 2 and Section 3 as well as Section4\_O and Section4\_M leaves the dataset as follows:

1		0	1	2	3	4	5	6	Section4_OM
2	0	30	125	1250	25	-10	505	228	1064
3	1	30	102	987	25	-12	505	228	1064
4	2	30	281	1543	50	-15	505	228	1064
5	3	30	101	1045	50	-9	505	228	1064
6	4	30	299	1754	100	-11	505	228	1064
7	5	30	134	961	100	-11	505	228	1064
8	10	45	201	2431	50	-9	1320	674	3048
9	11	45	201	2987	50	-11	1320	674	3048
10	12	45	451	2941	100	-11	1320	674	3048
11	13	45	451	1934	100	-9	1320	674	3048
12	14	45	876	1965	200	-11	1320	674	3048
13	19	125	149	2144	50	-15	10766	5400	24552
14	20	125	343	3871	100	-12	10766	5400	24552
15	21	125	343	2890	100	-9	10766	5400	24552
16	22	125	967	3813	200	-9	10766	5400	24552
17	23	125	967	2671	200	-8	10766	5400	24552

Figure 41. TiO<sub>2</sub> training dataset example by joining sections

About the second approach, consider Figure 37 and Figure 38. As an example, *Number O(6)* has been removed as it would interfere with the output expected. First four values from the NanoFingerprint have been also removed as explained in the previous approach. The remaining features after selecting most significant features in the example are as follows:

1	Name	Eng.Size (X0)	Size in Water (X1)	Size in PBS (X2)	Concentration (X4)	Zeta Potential (X5)	Y - LDH	Number Oxygen Atoms	Number Metal Atoms	Number O(1)	Number O(2)	Number O(3)	Number M(2)	Number M(3)	Number M(4)	Number M(5)	Number M(6)
2	TiO <sub>2</sub> _30	30	125	1250	25	-10	0.9	505	228	58	165	282	4	26	38	42	118
3	TiO <sub>2</sub> _30	30	102	987	25	-12	1	505	228	58	165	282	4	26	38	42	118
4	TiO <sub>2</sub> _30	30	281	1543	50	-15	0.75	505	228	58	165	282	4	26	38	42	118
5	TiO <sub>2</sub> _30	30	101	1045	50	-9	0.7	505	228	58	165	282	4	26	38	42	118
6	TiO <sub>2</sub> _30	30	299	1754	100	-11	1.04	505	228	58	165	282	4	26	38	42	118
7	TiO <sub>2</sub> _30	30	134	961	100	-11	1.09	505	228	58	165	282	4	26	38	42	118
8	TiO <sub>2</sub> _30	30	600	1876	200	-12	1.15	505	228	58	165	282	4	26	38	42	118
9	TiO <sub>2</sub> _30	30	298	1165	200	-12	1.2	505	228	58	165	282	4	26	38	42	118
10	TiO <sub>2</sub> _45	45	129	2567	25	-9	0.9	1320	674	60	368	892	0	96	108	128	342
11	TiO <sub>2</sub> _45	45	129	2309	25	-10	0.85	1320	674	60	368	892	0	96	108	128	342
12	TiO <sub>2</sub> _45	45	201	2431	50	-9	0.75	1320	674	60	368	892	0	96	108	128	342
13	TiO <sub>2</sub> _45	45	201	2987	50	-11	0.78	1320	674	60	368	892	0	96	108	128	342
14	TiO <sub>2</sub> _45	45	451	2941	100	-11	1.4	1320	674	60	368	892	0	96	108	128	342

Figure 42. TiO<sub>2</sub> dataset after preprocessing example

Applying the technique mentioned above across the whole subset, size is reduced to [24, 36]. Example of a training dataset considering split of 70:15:15, assigning 15% of samples to test, 15% to validation, and remaining to training subset has been added below:

1		0	1	2	3	4	5	6	7	8	9	16	17	18	19	20	24	25	26	122
2	12	45	451	2941	100	-11	1320	674	60	368	892	0	96	108	128	342	60	368	892	0
3	21	125	343	2890	100	-9	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
4	0	30	125	1250	25	-10	505	228	58	165	282	4	26	38	42	118	58	165	282	4
5	6	30	600	1876	200	-12	505	228	58	165	282	4	26	38	42	118	58	165	282	4
6	13	45	451	1934	100	-9	1320	674	60	368	892	0	96	108	128	342	60	368	892	0
7	22	125	967	3813	200	-9	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
8	20	125	343	3871	100	-12	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
9	18	125	149	3782	50	-10	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
10	19	125	149	2144	50	-15	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
11	23	125	967	2671	200	-8	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
12	9	45	129	2309	25	-10	1320	674	60	368	892	0	96	108	128	342	60	368	892	0
13	7	30	298	1165	200	-12	505	228	58	165	282	4	26	38	42	118	58	165	282	4
14	16	125	136	3215	25	-11	10766	5400	716	2820	7230	32	464	794	916	3194	716	2820	7230	32
15	2	30	281	1543	50	-15	505	228	58	165	282	4	26	38	42	118	58	165	282	4
16	11	45	201	2987	50	-11	1320	674	60	368	892	0	96	108	128	342	60	368	892	0
17	3	30	101	1045	50	-9	505	228	58	165	282	4	26	38	42	118	58	165	282	4

Figure 43. TiO2 training dataset example (Part 1)

1	122	131	140	149	158	13424	13433	13442	13487	13496	13505	13514	13523	13568	13577	13586	13595	13604
2	0	96	108	128	342	0	24	36	0	96	148	212	280	0	192	284	404	1372
3	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
4	4	26	38	42	118	2	28	28	4	32	64	80	150	4	46	86	102	438
5	4	26	38	42	118	2	28	28	4	32	64	80	150	4	46	86	102	438
6	0	96	108	128	342	0	24	36	0	96	148	212	280	0	192	284	404	1372
7	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
8	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
9	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
10	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
11	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
12	0	96	108	128	342	0	24	36	0	96	148	212	280	0	192	284	404	1372
13	4	26	38	42	118	2	28	28	4	32	64	80	150	4	46	86	102	438
14	32	464	794	916	3194	32	142	542	0	376	1074	1506	2684	64	1016	2070	2932	12114
15	4	26	38	42	118	2	28	28	4	32	64	80	150	4	46	86	102	438
16	0	96	108	128	342	0	24	36	0	96	148	212	280	0	192	284	404	1372
17	4	26	38	42	118	2	28	28	4	32	64	80	150	4	46	86	102	438

Figure 44. TiO2 training dataset example (Part 2)

#### 4.3.2. Papa – ZnO dataset

Same approaches are applied to the subset of only “ZnO” nanocompounds, in this case however, given the different composition of the nanocompound in comparison with “TiO2” dataset, the non-zero features differ. In this example, remaining subset size would be of [7, 147], leaving only 10 samples for validating the model, whereas when removing 0s the training subset would be of size [7, 11] when joining sections.

#### 4.3.3. Papa – TiO2+ZnO dataset

When considering all rows in the dataset, a possible configuration of available training samples when joining sections would be [27, 11]. Size when removing 0s is in this case of [27, 168].

#### 4.3.4. Anantha dataset

From Anatha\_db, no subset was required, in this case there were more nanocompounds to be considered, which implied that each model generated would remove different zeros depending on the available training samples used. To grasp the size of a possible available dataset, a possible training dataset observed was [306, 473], leaving a total of 54 samples for testing.

### 4.4. Model results

Main focus was put on multilinear regression as a quick and direct implementation was possible on the webportal since only the model coefficients had to be. Additional models have been investigated considering the definitions provided in Section 3.4. In this section the results obtained using only data-driven analysis as well as the combination of this analysis and recursive feature elimination are presented.

#### 4.4.1. Papa – TiO2 model

When using the joined dataset without performing any additional feature analysis, the following results were obtained:

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.9031	0.1161	0.0884	0.8234	0.5232	0.2256	0.1919	0.4578	0.5170	0.1997	0.1727	0.3159
1	0.8376	0.1475	0.1159	0.7205	0.9706	0.0796	0.0736	0.9348	0.3414	0.2195	0.2050	0.1729
2	0.8440	0.1344	0.1096	0.7301	0.8635	0.1768	0.1478	0.7399	0.6312	0.1686	0.1594	0.5122
3	0.8555	0.1236	0.0862	0.7474	0.7248	0.2335	0.1835	0.0646	0.3866	0.2497	0.1903	-0.0702
4	0.9246	0.0963	0.0595	0.8598	0.1987	0.3780	0.3350	0.0271	0.1979	0.2853	0.2176	-0.3966
5	0.9133	0.1162	0.0874	0.8405	0.2970	0.2442	0.1995	0.1722	0.4848	0.2277	0.1737	0.1104
6	0.8867	0.1227	0.0849	0.7964	0.7470	0.2509	0.2192	0.2413	0.2798	0.2522	0.2244	-0.0911
7	0.8073	0.1451	0.1176	0.6769	0.8505	0.1741	0.1549	0.4733	0.4763	0.2064	0.1843	0.2686
8	0.9005	0.1226	0.0888	0.8190	0.4364	0.2753	0.2473	-3.1821	0.3561	0.2096	0.1688	0.2461
9	0.9354	0.0895	0.0723	0.8786	0.4536	0.2655	0.2132	-1.2871	0.2280	0.2498	0.2193	-0.0709
10	0.8867	0.1227	0.0849	0.7964	0.7470	0.2509	0.2192	0.2413	0.2798	0.2522	0.2244	-0.0911
11	0.8441	0.1451	0.1123	0.7302	0.7272	0.1229	0.1162	0.4115	0.2543	0.2228	0.2097	0.1481
12	0.8826	0.1383	0.1029	0.7899	0.3316	0.1643	0.1587	-0.2011	0.5763	0.1986	0.1797	0.3231
13	0.8567	0.1259	0.0894	0.7493	0.8311	0.2003	0.1545	0.7468	0.4897	0.1972	0.1604	0.3324
14	0.8553	0.1420	0.1110	0.7471	0.9159	0.1247	0.1050	0.7444	0.5233	0.2070	0.1791	0.2645

Figure 45. Model metrics using TiO2 dataset joining sections

Analysing the r-squared metric for training, it is observed that nearly in all instances model represents with high fidelity the data provided in validation and test, however it fails to provide consistent results for the unseen data. When checking the r-squared for validation set, it is seen that in most cases model overfitted, since for those instances with negative or close to 0 r-squared, model fails to predict these values properly. Checking CCC, a similar pattern is observed. Only model 3 is seen to show consistency with a model not overfitted, as values for both CCC and r-squared do not extremely differ as seen with the remaining models. Considering MAE and RMSE, the difference across all predicted and expected samples does not show high discrepancy between expected and predicted value.

For the same model but including only the parameters that were flagged as the optimal for LDH prediction for TiO<sub>2</sub> (X0, X4), better results are obtained in comparison with the previous one. Models listed below show in most cases no overfit. Instances with negative r-squared for validation indicate that samples used in validation had no reference with the ones using in training, and model fails to fit the subset. Considering that samples were chosen at random, it was observed that validation samples belonged just to one of the nanocompounds, which lead to higher errors, as seen in both MAE and RMSE. Based on the features used, which were proven in [3] to provide the best results when using Linear Regression, models using only these parameters together with a compacted version of the fingerprint will mostly provide better predictions.

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.7925	0.1680	0.1300	0.6564	0.6847	0.1249	0.0911	0.5560	0.8308	0.1270	0.1199	0.7238
1	0.8191	0.1372	0.0997	0.6936	0.6901	0.2516	0.2008	0.5241	0.8170	0.1303	0.1229	0.7093
2	0.8401	0.1489	0.1145	0.7242	0.7548	0.2183	0.2087	0.4713	0.8144	0.1445	0.1264	0.6427
3	0.8729	0.1482	0.1068	0.7744	-0.3776	0.2198	0.1876	-1.1927	0.8617	0.1224	0.0975	0.7436
4	0.8103	0.1747	0.1465	0.6811	0.9272	0.0786	0.0702	0.7754	0.9081	0.0959	0.0942	0.8427
5	0.7637	0.1643	0.1296	0.6177	0.8519	0.1703	0.1313	0.6847	0.8455	0.1343	0.1177	0.6913
6	0.8624	0.1463	0.1182	0.7581	0.5742	0.2146	0.1393	0.0776	0.8979	0.1002	0.0976	0.8281
7	0.8498	0.1478	0.1043	0.7388	0.4986	0.2291	0.2014	-0.0159	0.8378	0.1363	0.1167	0.6821
8	0.8541	0.1601	0.1212	0.7454	-0.0067	0.1628	0.1390	-2.2051	0.8354	0.1314	0.1207	0.7045
9	0.8760	0.1409	0.1081	0.7793	0.0115	0.2344	0.1708	-0.8990	0.8381	0.1272	0.1208	0.7230
10	0.8875	0.1312	0.1022	0.7977	0.5466	0.2645	0.2394	-0.3257	0.9396	0.0768	0.0752	0.8991
11	0.8182	0.1575	0.1174	0.6924	0.7294	0.1761	0.1394	0.5264	0.8757	0.1069	0.1045	0.8044
12	0.7850	0.1641	0.1280	0.6461	0.8511	0.1399	0.0993	0.7554	0.8731	0.1128	0.1025	0.7824
13	0.7804	0.1745	0.1441	0.6399	0.9464	0.0755	0.0691	0.8871	0.8951	0.1014	0.0973	0.8241
14	0.8512	0.1460	0.1102	0.7409	0.4316	0.2069	0.1964	0.2575	0.9137	0.0922	0.0872	0.8545

Figure 46. TiO<sub>2</sub> modelling results using proposed features by authors

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.7834	0.1519	0.1269	0.6439	0.9466	0.0722	0.0565	0.9102	0.7358	0.2048	0.1797	0.6202
1	0.8618	0.1292	0.1053	0.7572	0.5858	0.1941	0.1642	0.2402	0.7367	0.2011	0.1678	0.6340
2	0.7589	0.1366	0.1125	0.6115	0.8485	0.1729	0.1504	0.7775	0.6657	0.2184	0.2002	0.5680
3	0.8332	0.1378	0.1062	0.7141	0.8212	0.1486	0.1378	0.5927	0.7704	0.1959	0.1720	0.6526
4	0.7833	0.1481	0.1216	0.6438	0.9090	0.1024	0.1011	0.8019	0.7375	0.2024	0.1734	0.6292
5	0.8472	0.1355	0.1008	0.7348	0.5817	0.2083	0.2063	0.0988	0.7695	0.1990	0.1763	0.6416
6	0.7977	0.1525	0.1314	0.6634	0.9227	0.0801	0.0666	0.8690	0.7124	0.2079	0.1835	0.6087
7	0.8444	0.1271	0.1063	0.7307	-0.1031	0.1889	0.1727	-0.7316	0.7229	0.2050	0.1769	0.6196
8	0.8066	0.1498	0.1207	0.6759	0.8606	0.0865	0.0805	0.7497	0.7452	0.2025	0.1748	0.6288
9	0.8444	0.1271	0.1063	0.7307	-0.1031	0.1889	0.1727	-0.7316	0.7229	0.2050	0.1769	0.6196
10	0.8472	0.1355	0.1008	0.7348	0.5817	0.2083	0.2063	0.0988	0.7695	0.1990	0.1763	0.6416
11	0.7834	0.1519	0.1269	0.6439	0.9466	0.0722	0.0565	0.9102	0.7358	0.2048	0.1797	0.6202
12	0.8644	0.1331	0.0976	0.7612	0.1022	0.2130	0.1807	-6.2376	0.7797	0.1954	0.1751	0.6542
13	0.8016	0.1504	0.1240	0.6689	0.8995	0.0958	0.0881	0.8093	0.7468	0.1994	0.1701	0.6399
14	0.8644	0.1331	0.0976	0.7612	0.1022	0.2130	0.1807	-6.2376	0.7797	0.1954	0.1751	0.6542

Figure 47. TiO<sub>2</sub> modelling results without features with std=0

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.8073	0.1386	0.1053	0.6768	0.7130	0.1669	0.1558	0.6371	0.6970	0.1800	0.1621	0.5040
1	0.8819	0.1327	0.1116	0.7887	0.1559	0.1929	0.1604	-0.2916	0.4595	0.2470	0.2227	0.0663
2	0.7092	0.1493	0.1236	0.5495	0.8928	0.1499	0.1412	0.8311	0.4980	0.2182	0.2063	0.2713
3	0.8047	0.1459	0.1210	0.6733	0.9047	0.1434	0.1374	0.8237	0.4981	0.2295	0.2148	0.1937
4	0.8779	0.1362	0.1080	0.7823	0.7327	0.1850	0.1805	0.1024	0.6279	0.2070	0.1832	0.3442
5	0.8540	0.1393	0.1119	0.7452	0.7889	0.1623	0.1539	0.6629	0.6580	0.1896	0.1732	0.4499
6	0.8227	0.1476	0.1159	0.6988	0.7341	0.1242	0.1220	0.3484	0.6105	0.2116	0.1882	0.3145
7	0.8434	0.1277	0.1145	0.7292	0.6239	0.2267	0.2006	0.5312	0.3353	0.2651	0.2453	-0.0759
8	0.8166	0.1405	0.1074	0.6901	0.8809	0.1591	0.1425	0.7992	0.6706	0.1872	0.1703	0.4634
9	0.8434	0.1277	0.1145	0.7292	0.6239	0.2267	0.2006	0.5312	0.3353	0.2651	0.2453	-0.0759
10	0.8318	0.1394	0.1085	0.7120	0.6553	0.1655	0.1486	0.3228	0.6548	0.1898	0.1746	0.4484
11	0.8600	0.1364	0.1038	0.7544	0.5338	0.1944	0.1758	-0.5882	0.6399	0.2059	0.1731	0.3513
12	0.8549	0.1452	0.1178	0.7465	0.8204	0.1444	0.1366	0.5771	0.5788	0.2168	0.1976	0.2807
13	0.8591	0.1503	0.1255	0.7529	0.7396	0.1262	0.1218	-0.2134	0.5494	0.2221	0.2046	0.2448
14	0.8587	0.1263	0.1132	0.7525	0.6614	0.2186	0.1879	0.5452	0.4214	0.2442	0.2223	0.0872

Figure 48. TiO<sub>2</sub> modelling results with proposed features in article

Performance of models when removing features with null standard deviation is more consistent across models. As per values of CCC and r-squared, models represent most of the samples in all instances. Even though at first glance it may appear that leaving only the proposed parameters for TiO<sub>2</sub> calculation + NanoFingerprint does perform better, results associated to test samples show that most models do not manage to represent this subset. R-squared is lower than 0.50 in all instances, but that does not necessarily imply model is overfitting, it just implies that it does not fully represent the test subset. The error linked to misclassification is higher as per RMSE and MAE. Overall, in comparison with using the previous approach, models do provide better model generalisation.

Main problems attributed to TiO<sub>2</sub> subset is linked to the small dataset available and high number of features. In most instances, especially when using joining sections approach, models tended to overfit and provide low results for both testing and validation.

As an alternative to the previous analysis, recursive feature elimination applied to the dataset was considered. recursively removing features that would not be useful for the model.

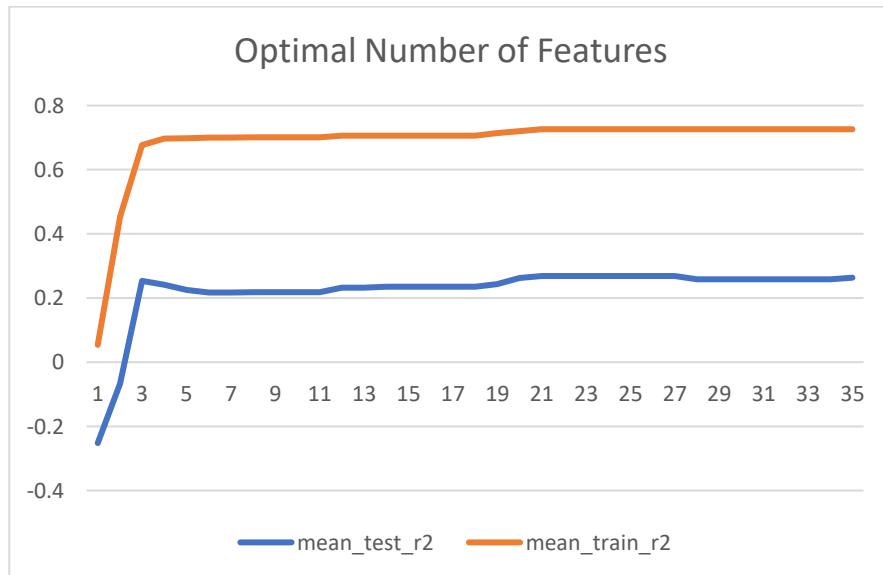


Figure 49. RFE applied to TiO<sub>2</sub> dataset

Considering the graph above, it is seen that the model performs better when using more than 20 features. Each number of features has a rank assigned, which helps to determine what model performed better and what number of features was used. In this case the best model was obtained when using 21 features. The metric used to validate these results was r-squared. It is seen that the second best in the ranking was when using 26 features. Below are the metrics associated to these 2 models:

#Features	CCC	RMSE	MAE	R <sup>2</sup>
21	0.8366444	0.175046	0.127087	0.593989
26	0.8366444	0.175046	0.127087	0.593989

Table 2. TiO2 results using RFE

Even though models appear to provide the same results for the metrics used, there's a slight difference. Full value is included in the appendix B.

#### 4.4.2. Papa – ZnO model

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	1	9.74E-16	8.66E-16	1	0.2136	0.1750	0.1685	-14.1301	0.8197	0.0512	0.0491	0.6367
1	1	3.91E-16	3.11E-16	1	0.1862	0.4569	0.3793	-68.0164	-0.7912	0.2411	0.2360	-7.0431
2	1	0	0	1	-0.2237	0.1236	0.1097	-0.6926	-0.0844	0.1842	0.1543	-3.6941
3	1	3.02E-16	2.89E-16	1	0.2136	0.1750	0.1685	-14.1301	0.8197	0.0512	0.0491	0.6367
4	1	1.79E-16	1.55E-16	1	-0.6667	0.1588	0.1426	-4.9669	-0.7698	0.2263	0.2160	-6.0867
5	1	4.97E-17	2.22E-17	1	0.1200	0.0979	0.0931	-3.7365	0.0583	0.1571	0.1374	-2.4147
6	1	1.72E-16	1.33E-16	1	-0.5303	0.4002	0.3983	-18.7740	0.5442	0.2004	0.1970	-4.5592
7	1	1.49E-16	1.11E-16	1	0.3251	0.1833	0.1673	-8.3335	0.6275	0.0892	0.0890	-0.1016
8	1	4.97E-17	2.22E-17	1	0.3251	0.1833	0.1673	-8.3335	0.6275	0.0892	0.0890	-0.1016
9	1	1.11E-16	6.66E-17	1	-0.0377	0.1210	0.1060	-584.2731	-0.7541	0.1866	0.1731	-3.8194
10	1	3.22E-16	2.22E-16	1	0.1862	0.4569	0.3793	-68.0164	-0.7912	0.2411	0.2360	-7.0431
11	1	0	0	1	0.6472	0.1204	0.1147	0.2049	-0.1990	0.2010	0.1579	-4.5925
12	1	1.11E-16	6.66E-17	1	0.2136	0.1750	0.1685	-14.1301	0.8197	0.0512	0.0491	0.6367
13	1	1.11E-16	6.66E-17	1	0.1200	0.0979	0.0931	-3.7365	0.0583	0.1571	0.1374	-2.4147
14	1	3.02E-16	2.89E-16	1	0.2136	0.1750	0.1685	-14.1301	0.8197	0.0512	0.0491	0.6367

Figure 50. ZnO modelling results using joined sections

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	1	5.73E-16	4.66E-16	1	-0.1110	0.3305	0.2950	-24.8494	0.0454	0.7532	0.7203	-225.9108
1	1	0	0	1	0.1266	0.1260	0.1211	0.0606	0.2596	0.0887	0.0858	-2.1449
2	1	1.11E-16	6.66E-17	1	-0.9204	0.1371	0.1363	-1.9375	-0.0110	0.2174	0.2103	-17.9055
3	1	1.72E-16	1.33E-16	1	-0.9204	0.1371	0.1363	-1.9375	-0.0110	0.2174	0.2103	-17.9055
4	1	1.11E-16	6.66E-17	1	0.5043	0.1584	0.1560	-1.5093	0.0287	0.2816	0.2803	-30.7199
5	1	9.93E-17	4.44E-17	1	-0.9204	0.1371	0.1363	-1.9375	-0.0110	0.2174	0.2103	-17.9055
6	1	7.02E-17	4.44E-17	1	0.6245	0.0981	0.0904	-2.1801	0.0277	0.2749	0.2735	-29.2370
7	1	1.99E-16	1.78E-16	1	0.4134	0.1841	0.1621	-10.2056	0.1929	0.2265	0.2145	-19.5179
8	1	4.28E-15	3.40E-15	1	0.0068	3.8248	3.5645	-11941.3584	0.0030	10.5325	10.0120	-44372.2126
9	1	2.33E-16	1.78E-16	1	0.2043	0.1704	0.1664	-71.6232	-0.2309	0.0857	0.0638	-1.9384
10	1	1.49E-16	1.11E-16	1	0.4134	0.1841	0.1621	-10.2056	0.1929	0.2265	0.2145	-19.5179
11	1	7.02E-17	4.44E-17	1	0.2043	0.1704	0.1664	-71.6232	-0.2309	0.0857	0.0638	-1.9384
12	1	4.68E-16	3.77E-16	1	0.2581	0.6572	0.6512	-42.1885	-0.0100	0.8172	0.8088	-266.0954
13	1	9.93E-17	4.44E-17	1	-0.1168	0.1288	0.1116	-7.1861	0.0255	0.2628	0.2609	-26.6339
14	1	4.28E-15	3.40E-15	1	0.0068	3.8248	3.5645	-11941.3584	0.0030	10.5325	10.0120	-44372.2126

Figure 51. ZnO modelling results using proposed features in article

As per the results seen above, model overfits in all instances, this is clearly seen analysing the results associated to the training dataset, where it is seen RMSE and MAE is close to 0. This implies that the models perfectly fit the data used for training. In contrast, results for both validation and test, especially when analysing the r-squared parameter show that these models do not represent these subsets. Negative r-squared is normally linked to missing bias, however all coefficients were considered when creating these models. The main reason behind these results is the low sample availability. Dataset is only formed by 9 samples, which makes it difficult to fit a model.

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	1	1.79E-16	1.11E-16	1	0.0017	0.2786	0.2777	-123.2258	-0.6332	0.3360	0.3331	-12.9364
1	1	1.07E-15	8.88E-16	1	0.2222	0.7083	0.6594	-54.5843	0.0532	0.6424	0.6412	-49.9444
2	1	3.18E-16	2.00E-16	1	0.1494	0.1205	0.1202	0.1404	-0.1216	0.0956	0.0955	-0.1273
3	1	6.30E-16	6.00E-16	1	-0.5938	0.2616	0.2611	-15.1993	0.4759	0.2538	0.2350	-6.9530
4	1	1.12E-15	9.77E-16	1	0.0062	0.2719	0.2688	-2956.7265	-0.3437	0.3724	0.2871	-16.1220
5	1	1.40E-16	8.88E-17	1	0.3716	0.1443	0.1391	-1.0817	-0.8807	0.2298	0.2273	-5.5176
6	1	1.99E-16	1.78E-16	1	0.2315	0.1504	0.1504	-55.5834	-0.4788	0.1130	0.1130	-0.5755
7	1	6.28E-16	5.33E-16	1	-0.6026	0.2835	0.2689	-11.5616	0.1991	0.1423	0.1280	-1.4988
8	1	9.93E-17	4.44E-17	1	0.1286	0.1218	0.1216	0.1221	-0.3198	0.1049	0.1048	-0.3581
9	1	5.62E-15	5.06E-15	1	-0.0255	0.7656	0.7438	-161.8103	0.2203	0.7166	0.7162	-62.3939
10	1	1.58E-15	1.38E-15	1	-0.5584	0.2719	0.2687	-16.4948	-0.4573	0.4490	0.4415	-23.8877
11	1	1.79E-15	1.55E-15	1	0.2222	0.7083	0.6594	-54.5843	0.0532	0.6424	0.6412	-49.9444
12	1	1.72E-16	1.33E-16	1	-0.9647	0.0895	0.0887	-2.9559	-0.8819	0.2294	0.2270	-5.4981
13	1	5.62E-15	5.06E-15	1	-0.0255	0.7656	0.7438	-161.8103	0.2203	0.7166	0.7162	-62.3939
14	1	1.07E-15	8.88E-16	1	0.2222	0.7083	0.6594	-54.5843	0.0532	0.6424	0.6412	-49.9444

Figure 52. ZnO modelling results without features with  $std = 0$

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	1	2.31E-15	2.22E-15	1	-0.2261	0.3911	0.3114	-49.5685	0.0074	0.3983	0.3817	-6343.6067
1	1	1.07E-15	1.04E-15	1	-0.4669	0.1668	0.1615	-21.7050	-0.0671	0.1263	0.1054	-637.3584
2	1	2.67E-16	2.00E-16	1	-0.7460	0.1851	0.1711	-3.7421	-0.0420	0.1192	0.1016	-567.4191
3	1	8.31E-16	6.22E-16	1	-0.2506	0.3596	0.2768	-11.9291	0.0080	0.2987	0.2912	-3568.6766
4	1	3.48E-16	2.44E-16	1	-0.0762	0.1869	0.1728	-37.8077	-0.0307	0.0515	0.0499	-105.0978
5	1	1.70E-15	1.47E-15	1	-0.2285	0.4036	0.3006	-24.4573	0.0062	0.3448	0.3377	-4754.9200
6	1	3.88E-16	2.89E-16	1	-0.0307	0.2634	0.2554	-55.6288	-0.0428	0.0956	0.0856	-364.6176
7	1	5.41E-15	4.15E-15	1	0.0011	0.2048	0.2044	-418.3495	0.0135	0.3822	0.3285	-5841.8663
8	1	6.90E-16	5.55E-16	1	-0.0167	0.2466	0.2392	-607.2517	-0.0307	0.1736	0.1455	-1203.9182
9	1	7.51E-16	6.00E-16	1	-0.0744	0.1954	0.1808	-41.4316	-0.0601	0.0459	0.0427	-83.3262
10	1	2.81E-16	1.78E-16	1	-0.0727	0.2432	0.2193	-93.6542	-0.0996	0.1047	0.1042	-437.1006
11	1	5.25E-16	4.44E-16	1	-0.2506	0.3596	0.2768	-11.9291	0.0080	0.2987	0.2912	-3568.6766
12	1	2.03E-15	1.47E-15	1	-0.2285	0.4036	0.3006	-24.4573	0.0062	0.3448	0.3377	-4754.9200
13	1	8.31E-16	6.22E-16	1	-0.2506	0.3596	0.2768	-11.9291	0.0080	0.2987	0.2912	-3568.6766
14	1	3.88E-16	3.77E-16	1	-0.0762	0.1869	0.1728	-37.8077	-0.0307	0.0515	0.0499	-105.0978

Figure 53. ZnO modelling results using proposed features in article

Similar to the previous approach, models overfit – only 2 features have been removed from previous instance. Distance between predicted samples and actual values are seen to be considerable as per RMSE and MAE.

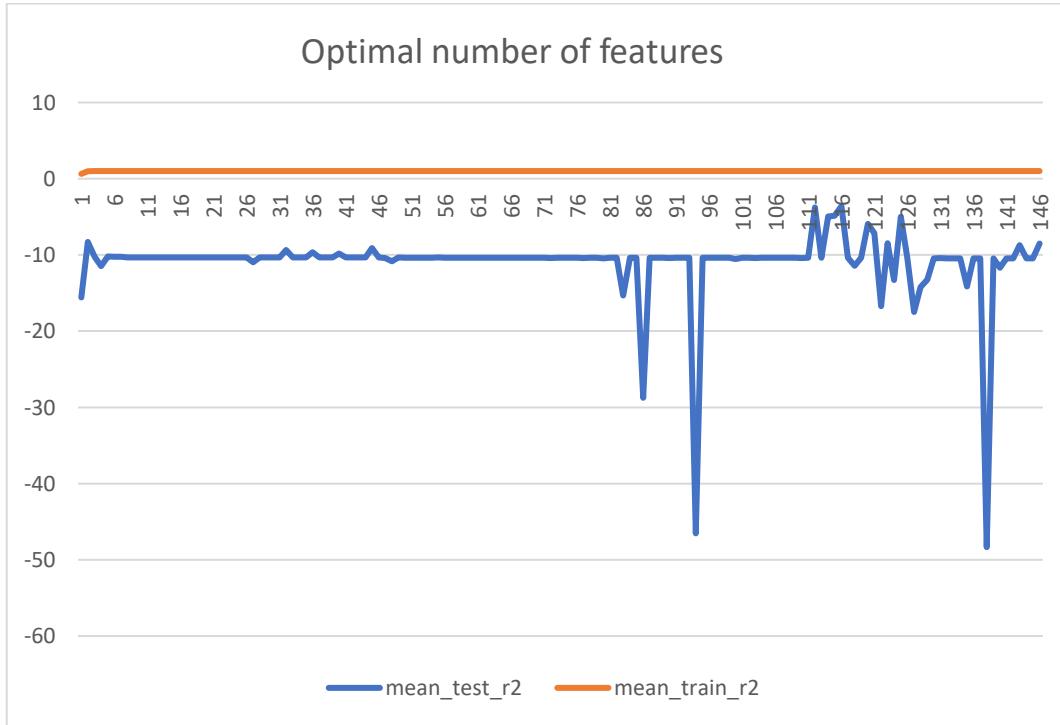


Figure 54. RFE applied to ZnO dataset

Considering the graph above, we observe similar results as the ones seen using data-driven analysis. Low sample availability for subset does not provide reliable models. Notice the spikes when using either 86, 94 or 138 features. Based on the metrics gathered from RFE, it is seen that 2 models had a better performance than the rest, however there are still not fitting the data analysed.

#Features	CCC	RMSE	MAE	R <sup>2</sup>
112	0.43203	0.40218	0.34431	-4.94128
116	0.43203	0.40218	0.34431	-4.94128

5. Table 3. ZnO results using RFE

#### 4.4.3. Papa – TiO2+ZnO model

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.8468	0.1388	0.1104	0.7343	0.7598	0.1259	0.1228	0.4435	0.5001	0.1518	0.1511	-0.2623
1	0.8817	0.1129	0.0833	0.7885	0.4220	0.2471	0.2010	-0.4748	0.4686	0.1788	0.1775	-0.7506
2	0.8921	0.1157	0.0870	0.8052	0.3669	0.2447	0.1837	-0.3141	0.4949	0.1675	0.1650	-0.5366
3	0.8439	0.1264	0.1016	0.7300	0.8097	0.1804	0.1540	0.6291	0.5944	0.1327	0.1223	0.0354
4	0.8357	0.1391	0.1040	0.7178	0.7350	0.1486	0.1181	0.4524	0.5754	0.1340	0.1317	0.0165
5	0.8704	0.1281	0.1021	0.7706	0.4978	0.1750	0.1349	0.0573	0.6007	0.1391	0.1341	-0.0600
6	0.8601	0.1267	0.0844	0.7546	0.5200	0.2072	0.1819	-0.4715	0.4529	0.1641	0.1618	-0.4743
7	0.8535	0.1302	0.0990	0.7444	0.3878	0.1842	0.1457	-1.1609	0.5336	0.1510	0.1479	-0.2481
8	0.8061	0.1377	0.1159	0.6751	0.9149	0.1225	0.1003	0.8328	0.5490	0.1476	0.1411	-0.1931
9	0.8538	0.1384	0.1094	0.7448	0.6385	0.1641	0.1568	-0.6076	0.6468	0.1377	0.1273	-0.0389
10	0.9383	0.0881	0.0785	0.8838	-0.2190	0.3004	0.2407	-1.2416	0.5675	0.1406	0.1349	-0.0828
11	0.8850	0.1130	0.0873	0.7938	0.3201	0.2361	0.2236	-0.0663	0.7051	0.1020	0.0955	0.4305
12	0.8203	0.1385	0.1114	0.6954	0.7890	0.1426	0.1413	0.4441	0.6032	0.1410	0.1361	-0.0891
13	0.8468	0.1388	0.1104	0.7343	0.7598	0.1259	0.1228	0.4435	0.5001	0.1518	0.1511	-0.2623
14	0.8606	0.1372	0.1027	0.7553	0.2468	0.1708	0.1453	-3.7201	0.5642	0.1467	0.1401	-0.1788

Figure 55. TiO2+ZnO modelling results when joining sections

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.8132	0.1442	0.1173	0.6853	0.8072	0.1287	0.0975	0.6737	0.7779	0.1150	0.1016	0.6487
1	0.7944	0.1450	0.1197	0.6589	0.8993	0.1115	0.0759	0.8317	0.7511	0.1227	0.1092	0.6003
2	0.7791	0.1402	0.1099	0.6381	0.8849	0.1414	0.1055	0.8029	0.6616	0.1509	0.1389	0.3956
3	0.8256	0.1237	0.0975	0.7176	0.6969	0.2104	0.1537	0.5499	0.7206	0.1271	0.1117	0.5710
4	0.8520	0.1266	0.0928	0.7421	0.5507	0.2048	0.1483	0.4087	0.6810	0.1439	0.1284	0.4502
5	0.8225	0.1499	0.1239	0.6985	0.8942	0.0663	0.0540	0.6587	0.7719	0.1226	0.1137	0.6009
6	0.7874	0.1498	0.1210	0.6493	0.9633	0.0654	0.0569	0.9326	0.7500	0.1249	0.1110	0.5856
7	0.8454	0.1338	0.0964	0.7322	0.5175	0.1863	0.1751	-0.2266	0.6161	0.1755	0.1732	0.1822
8	0.8320	0.1441	0.1071	0.7124	0.1391	0.1315	0.1048	-3.5341	0.7739	0.1224	0.1137	0.6021
9	0.8666	0.1136	0.0804	0.7646	0.4863	0.2293	0.1975	0.3562	0.7125	0.1321	0.1131	0.5364
10	0.8218	0.1324	0.0996	0.6976	0.6259	0.1786	0.1322	0.0469	0.6788	0.1522	0.1417	0.3847
11	0.8516	0.1271	0.0944	0.7416	0.5667	0.2018	0.1370	0.4140	0.6883	0.1421	0.1269	0.4634
12	0.8165	0.1309	0.0998	0.6899	0.8456	0.1720	0.1372	0.7105	0.7741	0.1202	0.1083	0.6160
13	0.8265	0.1438	0.1082	0.7043	0.7873	0.1278	0.1131	0.4984	0.7918	0.1205	0.1165	0.6142
14	0.8265	0.1438	0.1082	0.7043	0.7873	0.1278	0.1131	0.4984	0.7918	0.1205	0.1165	0.6142

Figure 56. TiO2+ZnO modelling results using proposed features in article

When joining both dataset under a single set, it is observed that no models fit the analysed data, considering CCC and r-squared metrics, it is seen that the results associated to training are way higher than the other sets. In both scenarios, RMSE and MAE should low impact on error. Regarding the second model, better results are obtained, in this case we are using the features proposed by the paper as well as a compacted NanoFingerprint. There are a couple of scenarios where r-squared is seen to be negative, which again indicates overfitting.

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.8669	0.1256	0.0943	0.7651	0.4588	0.1926	0.1548	-0.2908	0.6704	0.1713	0.1447	0.5038
1	0.8234	0.1408	0.1055	0.6998	0.3075	0.1198	0.1075	-1.2067	0.8039	0.1316	0.1169	0.7071
2	0.8156	0.1233	0.0946	0.6887	0.7559	0.1953	0.1837	0.6276	0.6369	0.1754	0.1512	0.4797
3	0.8202	0.1414	0.1077	0.6951	0.8523	0.1156	0.0948	0.6243	0.8237	0.1263	0.1141	0.7304
4	0.7609	0.1386	0.1071	0.6140	0.9078	0.1265	0.1070	0.8514	0.6871	0.1601	0.1368	0.5667
5	0.8392	0.1382	0.1072	0.7229	0.6346	0.1488	0.1418	0.0351	0.7514	0.1495	0.1132	0.6220
6	0.8534	0.1280	0.1075	0.7442	0.6274	0.1901	0.1704	0.0282	0.8195	0.1280	0.1131	0.7229
7	0.8278	0.1380	0.1088	0.7062	0.5748	0.1626	0.1568	0.1206	0.7046	0.1625	0.1195	0.5533
8	0.8392	0.1382	0.1072	0.7229	0.6346	0.1488	0.1418	0.0351	0.7514	0.1495	0.1132	0.6220
9	0.7609	0.1386	0.1071	0.6140	0.9078	0.1265	0.1070	0.8514	0.6871	0.1601	0.1368	0.5667
10	0.7884	0.1398	0.1089	0.6507	0.8604	0.1377	0.1306	0.7294	0.6084	0.1842	0.1539	0.4262
11	0.7609	0.1386	0.1071	0.6140	0.9078	0.1265	0.1070	0.8514	0.6871	0.1601	0.1368	0.5667
12	0.8499	0.1324	0.1010	0.7390	0.4084	0.1715	0.1637	-0.5640	0.7632	0.1483	0.1326	0.6279
13	0.8336	0.1308	0.1090	0.7144	-1.85E-12	40598895679	18156378110	-2.61E+22	-7.49E-13	57415508886	36312756220	-5.57E+22
14	0.8597	0.1252	0.1042	0.7549	0.4323	0.1985	0.1869	-0.2548	0.7854	0.1419	0.1185	0.6595

Figure 57. TiO<sub>2</sub>+ZnO modelling results without features with std = 0

Checking the model below, where those features with std equal to 0 have been removed, it is observed that multiple models have good CCC and r-squared metrics. RMSE and Mae is higher than previous models, but models still manage to fit most of the unseen data.

Index	train_CCC	train_RMSE	train_MAE	train_r2_score	validation_CCC	validation_RMSE	validation_MAE	validation_r2_score	test_CCC	test_RMSE	test_MAE	test_r2
0	0.8390	0.1210	0.0966	0.7227	0.8824	0.1355	0.1285	0.6974	0.6012	0.2719	0.2638	-0.2532
1	0.9097	0.0949	0.0843	0.8344	0.3812	0.2203	0.1677	-0.5622	0.6071	0.2760	0.2667	-0.2910
2	0.8545	0.1110	0.0922	0.7460	0.8091	0.1628	0.1486	0.5503	0.6154	0.2752	0.2695	-0.2840
3	0.8892	0.1088	0.0890	0.7998	0.2086	0.1722	0.1362	-0.4319	0.5597	0.2842	0.2765	-0.3694
4	0.8766	0.1020	0.0798	0.7803	0.6677	0.1896	0.1721	0.5771	0.6635	0.2218	0.2117	0.1661
5	0.9019	0.0964	0.0840	0.8214	0.4732	0.2152	0.1569	0.0798	0.6156	0.2628	0.2512	-0.1707
6	0.9130	0.0943	0.0779	0.8399	0.5162	0.2358	0.2049	-0.3281	0.4870	0.3281	0.3175	-0.8250
7	0.8754	0.1088	0.0909	0.7784	0.4380	0.1743	0.1303	0.0687	0.5613	0.2802	0.2702	-0.3304
8	0.8912	0.1043	0.0877	0.8037	0.6283	0.2048	0.1699	-0.0145	0.5184	0.3219	0.3123	-0.7563
9	0.8545	0.1110	0.0922	0.7460	0.8091	0.1628	0.1486	0.5503	0.6154	0.2752	0.2695	-0.2840
10	0.8826	0.0996	0.0784	0.7899	0.4869	0.1818	0.1619	0.3913	0.5733	0.2677	0.2584	-0.2142
11	0.8690	0.1163	0.0943	0.7683	0.8079	0.1346	0.1091	0.4885	0.5876	0.2779	0.2707	-0.3086
12	0.8373	0.1237	0.1020	0.7202	0.9400	0.0922	0.0747	0.8344	0.6153	0.2679	0.2601	-0.2160
13	0.8657	0.1221	0.1004	0.7632	0.6234	0.0967	0.0810	0.3193	0.5504	0.2935	0.2859	-0.4596
14	0.9235	0.0849	0.0695	0.8579	0.3409	0.2616	0.2192	0.0378	0.7042	0.2003	0.1653	0.3203

Figure 58. TiO<sub>2</sub>+ZnO modelling results when using proposed features in article

When using the parameters advised within the paper as optimal for regression models, it is seen that r-squared for testing responds to negative values, which indicates that models create do not fit the data analysed. Checking CCC and r-squared parameters for training, we can observe a huge difference between values when comparing it to testing samples. This is an indicator of possible model being overfitted.

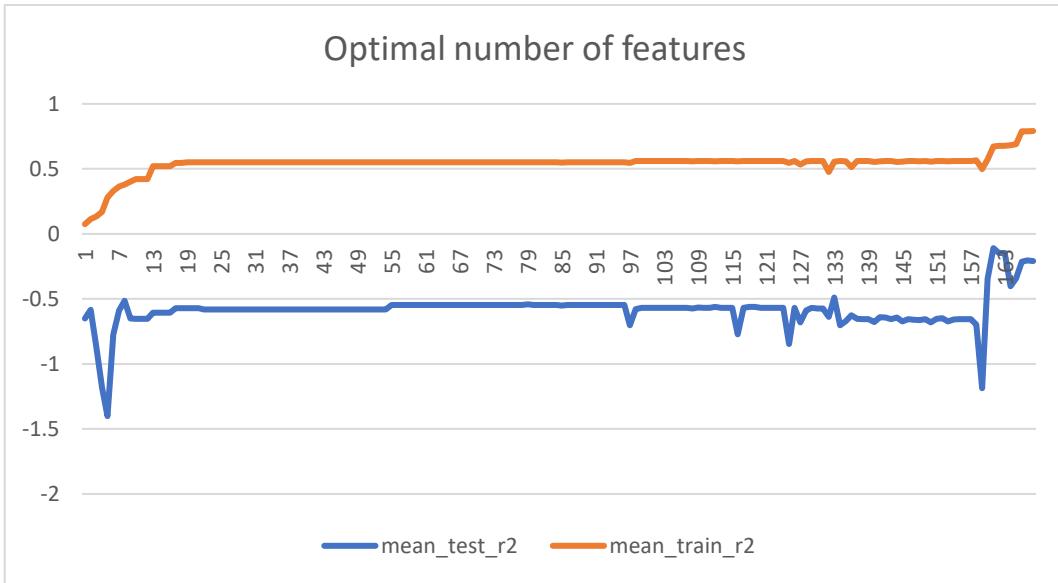


Figure 59. RFE applied to  $TiO_2+ZnO$  dataset

When applying RFE to the whole dataset, we observe that all r-squared values for the validation dataset is below 0. No changes are perceived in model performance between using 13 features and using 150 approximately. As part rank provided by the technique, it was observed that best ranked features were found when using 160 or more feature. In this case we had around 168 features to analyse and find the optimal ones for the model. When checking the results obtained for the case of use 161 and 163, we observed the following:

#Features	CCC	RMSE	MAE	$R^2$
161	0.3986080	0.2526117	0.20380952	0.23287
163	0.2544713	0.3305967	0.2968154	-0.313

Table 4.  $TiO_2+ZnO$  results using RFE

#### 4.4.4. Anantha model

Index	train_RMSE	train_MAE	train_r2_score	validation_RMSE	validation_MAE	validation_r2_score	test_RMSE	test_MAE	test_r2
0	0.40336	0.16270	-0.80543	0.30429	0.09259	-0.22921	0.50452	0.25455	-1.52393
1	0.48795	0.23810	-1.16708	0.47140	0.22222	-1.31423	0.57208	0.32727	-1.41279
2	0.33923	0.11508	0.31535	0.36004	0.12963	-0.14136	0.38139	0.14545	0.18161
3	0.35074	0.12302	0.10261	0.38490	0.14815	-0.06764	0.40452	0.16364	-0.06348
4	0.50787	0.25794	-3.55007	0.49065	0.24074	-2.36793	0.50452	0.25455	-3.19601
5	0.30861	0.09524	0.34585	0.33333	0.11111	-0.01112	0.40452	0.16364	0.16322
6	0.38832	0.15079	-0.32849	0.40825	0.16667	-0.03909	0.42640	0.18182	-0.21833
7	0.32121	0.10317	0.34251	0.36004	0.12963	0.21424	0.40452	0.16364	0.20704
8	0.35635	0.12698	-0.05856	0.40825	0.16667	-1.15287	0.44721	0.20000	-0.09383
9	0.30211	0.09127	0.32460	0.40825	0.16667	0.25139	0.40452	0.16364	0.18478
10	0.29547	0.08730	0.33676	0.36004	0.12963	0.11367	0.38139	0.14545	0.15736
11	0.33923	0.11508	0.31535	0.36004	0.12963	-0.14136	0.38139	0.14545	0.18161
12	0.38832	0.15079	-0.32849	0.40825	0.16667	-0.03909	0.42640	0.18182	-0.21833
13	0.32733	0.10714	0.31940	0.36004	0.12963	-0.49289	0.42640	0.18182	-0.07571
14	0.38318	0.14683	-0.42213	0.43033	0.18519	-0.65147	0.42640	0.18182	-0.41614

Figure 60. Anantha modelling results

Given that Anantha is formed by multiple nanocompounds, using LinearRegression does not manage to fully fit all sets. Considering r-squared, it is seen that most of the model provide negative value for this metric. In this case, metric CCC has not been used and analysis of models has been performed considering the other 3 remaining metrics. The decision behind this exclusion is linked to the fact that the toxicity prediction associated with this article is categorical – if the output has a value higher than 0.5 then nanocompound is considered toxic.

After presenting all models below is a table with metrics comparing our results with the results present in the different articles:

<b>Study</b>	<b>CCC</b>	<b>RMSE</b>	<b>r-squared</b>
Papa – TiO2 LinearRegression	0.94	0.10	0.90
Our Study: Nanofingerprint + Papa  (Joined Sections)	0.5763	0.1986	0.3231
Our Study: Nanofingerprint + Papa  (Proposed parameters + Joined Sections)	0.9396	0.07	0.8991
Our Study: Nanofingerprint + Papa  (Feature removal with std=0)	0.7797	0.1954	0.6542
Our Study: Nanofingerprint + Papa  (Proposed parameters + feature removal)	0.697	0.18	0.50
Our Study: NanoFingerprint + Papa (RFE)	0.8364	0.175	0.5939

*Figure 61. Compiled results for TiO2*

<b>Study</b>	<b>CCC</b>	<b>RMSE</b>	<b>r-squared</b>
Papa – ZnO LinearRegression	0.99	0.04	0.97
Our Study: Nanofingerprint + Papa  (Joined Sections)	0.8197	0.0512	0.6367
Our Study: Nanofingerprint + Papa	0.25	0.08	-2.1

(Proposed parameters + Joined Sections)			
Our Study: Nanofingerprint + Papa  (Feature removal with std=0)	-0.3198	0.10	-0.3585
Our Study: Nanofingerprint + Papa  (Proposed parameters + feature removal)	-0.06	0.04	-83
Our Study: NanoFingerprint + Papa (RFE)	0.43	0.40	-4.94

Figure 62. Compiled results for ZnO

Study	CCC	RMSE	r-squared
Papa – TiO <sub>2</sub> +ZnO LinearRegression	0.92	0.08	0.96
Our Study: Nanofingerprint + Papa  (Joined Sections)	0.7051	0.1020	0.43
Our Study: Nanofingerprint + Papa  (Proposed parameters + Joined Sections)	0.7918	0.12	0.61
Our Study: Nanofingerprint + Papa  (Feature removal with std=0)	0.8237	0.1263	0.7304
Our Study: Nanofingerprint + Papa  (Proposed parameters + feature removal)	0.7042	0.2003	0.3203
Our Study: NanoFingerprint + Papa (RFE)	0.39	0.25	0.23

Figure 63. Compiled results for TiO<sub>2</sub>+ZnO

## 6. Conclusions & Future Work

In this project multiple functionalities have been developed for the production of the web portal framework. An initial assessment was performed to understand the required for the webportal development. Website is currently up and running and publicly accessible enabling users to compute NanoFingerprints or determine the Toxicity level for a given nanocompound. Back-end has been developed considering any future modifications in terms of adding either new functionalities or providing additional files to the end user. By creating subdirectories for each interaction with the portal itself we have ensured that the same data is not accessed by different users at the same time. The development of a local environment allowed to test any functionality without compromising access to the web portal.

The analysis of different toxicity prediction algorithms has been performed and further research on the usage of the NanoFingerprint for regression models have been presented. As per results obtained it is remarkable to state that the limitation on samples available for training and testing the model has led to poor results in comparison to the ones provided in the multiple articles studied. The introduction of Recursive Feature Elimination as means of selecting features that could benefit the model has provided a different perspective. It was observed for example that on TiO<sub>2</sub>+ZnO dataset the number of features with the highest rank was around 160 features. When checking the features that were selected, it was observed that most of these were located in the Section 4 of the NanoFingerprint. Further analysis with models mostly considering Section 4 could be beneficial in testing further models for toxicity prediction. Additionally, methods based on feature selection by means of ranking analysis can also provide a different insight on what part of the NanoFingerprint do provide further relevance to models. Due to time constraints, no further models could be analysed, however an extension of the models and focus on the utilization of algorithms based on Support Vector Machines would be an alternative to Linear Regression, which could at the same time be added to the website as a complement to the current model in place.

In regard to the website, future work should be focused on the visualisation of the shell as a mean to verify the results obtained from the NanoFingerprint. Multiple works are available using JavaScript and the integration of such functionality would not impact current resources available in the server.

## 7. References

- [1] F. Serratosa, S. Álvarez, L. Escorihuela and M. Calatayud, "Subgraph NanoFingerprint for modelling metal oxide nanoparticles based on connected atoms exploration," 2022.
- [2] A. Nel, T. Xia, L. Mädler and N. Li, "Toxic potential of materials at the nanolevel," *Science*, pp. 622-627, 2006.
- [3] E. Papa, J. P. Doucet and A. Doucet-Panaye, "Linear and non-linear modelling of the cytotoxicity of TiO<sub>2</sub> and ZnO nanoparticles by empirical descriptors," *SAR and QSAR in Environmental Research*, vol. 26, pp. 647-665, 2015.
- [4] N. Anantha and A. Palaniappan, "NanoTox: Developemnt of a Parsimonious In Silico Model for Toxicity Assessment of Metal-Oxide Nanoparticles Using PhysicoChemical Features," *ACS Omega*, pp. 11729-11739, 2021.
- [5] A. Gajewicz, N. Schaeublin, B. Rasulev, S. Hussain, D. Leszczynska, T. Puzyn and J. Leszczynski, "Towards understanding mechanisms governing cytotoxicity of metal oxides nanoparticles: Hints from nano-QSAR studies," *Nanotoxicology*, pp. 313-325, 2015.
- [6] T. Puzyn, B. Rasulev, A. Gajewicz, X. Hu, T. Dasari and A. Michalkova, "Using nano-QSAR to predict the cytotoxicity of metal oxide nanoparticles.,," *Nat Nanotechnol*, vol. 6, 2011.
- [7] K. Tämm, L. Sikk, J. Burk, R. Rallo, S. Pokhrel, L. Mädler, J. Scott-Fordsmund and P. Burk, "Parametrization of nanoparticles: development of full-particle nanodescriptors," *Nanoscale*, 2016.
- [8] F. Effendy and B. Adhilaksono, "Performance Comparison of Web Backend and Database: A Case Study of Node.JS, Golang and MySQL, Mongo DB," *Recent Advances in Computer Science and Communications*, vol. 14, no. 6, pp. 1955-1961, 2021.
- [9] I. Guyon and A. Elisseeff, "An Introduction of Variable and Feature Selection," *Journal of Machine Learning Research*, vol. 3, pp. 1157-1182, 2003.
- [10] N. Chirico and P. Gramatica, "Real External Predictivity of QSAR Models: How To Evaluate It? Comparison of Different Validation Criteria and Proposal of Using the Concordance Correlation Coefficient," *Journal of Chemical Information and Modelling*, pp. 2320-2335, 2011.
- [11] C. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [12] R. Duda, P. Hart and D. Stork, *Pattern Classification*, Second ed., John Wiley & Sons, Inc, 2001.
- [13] G. Varoquaux, "Cross-validation failure: small sample sizes lead to large error bars," 2017.

- [14] “Golang Documentation,” [Online]. Available: <https://go.dev/doc/>. [Accessed 3 March 2023].
- [15] “Shell Depth Calculator,” [Online]. Available: <https://nanogen.me/shell-depth>. [Accessed 7 April 2023].
- [16] I. Guyon, J. Weston, S. Barnhill and V. Vapnik, “Gene Selection for Cancer Classification using Support Vector Machines,” *Machine Learning*, vol. 46, pp. 389-422, 2002.

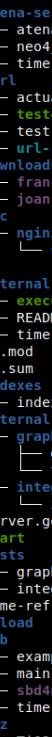
## Appendix A

## Access to Atena and Aura

Both Atena and Aura are Linux hosts located in the university labs. As stated in previous sections, Atena is the device in charge of hosting the website. As a security measure, any changes performed in Atena in terms of the website have to be performed through Aura.

In the scenario of requiring modifications in the workflow or presentation of data for the website, file are required to be uploaded in “atenaserver” directory, located in Aura.

Actividades 14 de jun 16:28



```
• Configurando tree (1.8.0-1) ...
Procesando disparadores para man-db (2.9.1-1) ...
user@user-HP-Compaq-6005-Pro-MT-PC:~$ cerrar sesión
atena@user-HP-Compaq-6005-Pro-MT-PC:~/atenaserver$ tree .
.
├── atena-services
│   ├── atenaserver.service
│   ├── neo4j.service
│   └── time
├── curl
│   ├── actual-url
│   │   └── test4services.sh
│   ├── test-url
│   │   └── url-rest -> actual-url
│   └── downloads
│       ├── francesc.png
│       └── joan.png
└── etc
    ├── nginx
    │   └── sites-available
    │       └── ckan
    ├── external
    │   ├── executor
    │   ├── README
    │   └── time
    ├── go.mod
    ├── go.sum
    └── indexes
        └── indexes.cql
    └── internal
        ├── graphdb
        │   ├── graphdb.go
        │   └── time
        └── integrity
            └── integrity.go
    └── server.go
    └── start
    └── tests
        ├── graphdb.go
        └── integrity.go
    └── time-ref
    └── upload
    └── web
        ├── example.grf
        ├── main.html
        ├── sbd4nano.png
        └── time
    └── xyz
        ├── TiO2_007.xyz
        ├── TiO2_010.xyz
        ├── TiO2_015.xyz
        └── ZnO_3nm.xyz
15 directories, 32 files
atena@user-HP-Compaq-6005-Pro-MT-PC:~/atenaserver$
```

Changes require of username and password, which can be provided by the local site admin providing a justification, as site is publicly accessible, and any unauthorized changes can lead to security violations.

## Web Portal: GenVector function

```

func GenVector(file string, shellThickness float64, genXYZ bool, maxBounds int, path string, strId string, base_name string) (fileRet string, fileRetNv string, err error, fileShell string) {
    defer func() { ...
    }()
    shellThickness = shellThickness * 10
    var xyz []Atom
    content, err := os.ReadFile(file)
    if err != nil { ...
    }
    lines := strings.Split(string(content), "\n")
    // check first line of .xyz file
    var head string
    if len(strings.Fields(lines[1])) > 3 { ...
    } else { ...
    }
    xyz = make([]Atom, len(lines))
    var median Atom

    for i, l := range lines { ...
    }
    numXYZ := float64(len(xyz))
    median.X /= numXYZ
    median.Y /= numXYZ
    median.Z /= numXYZ
    maxRadius := maxRadiusDistInit(median, xyz)
    shell := make(map[int]Atom)
    shellLine := maxRadius - shellThickness
    if shellLine <= 0.0 { ...
    }
    nearShellLine := shellLine - min_distance
    if nearShellLine <= 0.0 { ...
    }
    length_shell := 0
    for i, _ := range xyz { ...
    }

    visited := make(map[string]bool)
    bounds := make([]Bound, 0, 100) // bound format IDATOM-IDATOM. Ex: 22-66
    for id, atom := range shell { ...
    }
    fmt.Println("Atoms ", len(shell))
    fmt.Println("bounds: ", len(bounds))

    NF_path := filepath.Join(path, strId)
    if err_dir := os.MkdirAll(NF_path, os.ModePerm); err_dir != nil {
        log.Fatal(err)
    }
    trimmed_basename := strings.TrimSuffix(base_name, filepath.Ext(base_name))
    fileZIP := genXYZ(shell, head, length_shell, NF_path, strId, trimmed_basename)
    fileToDownload, fileToDownloadNoVer, err := computeShell(shell, bounds, driver, shellThickness, maxBounds, maxRadius*2, NF_path, strId, trimmed_basename)
    return fileToDownload, fileToDownloadNoVer, err, fileZIP
}

func computeShell map[int]Atom, bounds []Bound, driver neo4j.Driver, shellThickness float64, maxBounds int, size float64, vectorPath string, strId string, fileName string) (string, string, error) {
    fileOut := fileName + " " + fmt.Sprintf("%v", shellThickness) + "nm_" + fmt.Sprintf("%v", maxBounds) + "_Verbose_NanoFingerprint.txt" // verbose fingerprint file - change from filename_base[2] from regex to filename
    fileOutNoVerbose := fileName + " " + fmt.Sprintf("%v", shellThickness) + "nm_" + fmt.Sprintf("%v", maxBounds) + ".NanoFingerprint.txt" // no verbose fingerprint file
    path_NF_Verbose := filepath.Join(vectorPath, fileOut)
    f, err := os.Create(path_NF_Verbose)
    path_NF := filepath.Join(vectorPath, fileOutNoVerbose)
    f2, err := os.Create(path_NF)

    check(err)
    defer f.Close()
    defer f2.Close()
    w := bufio.NewWriter(f)
    w2 := bufio.NewWriter(f2)

    sessionW := driver.NewSession(neo4j.SessionConfig{AccessMode: neo4j.AccessModeWrite})
    defer unsafeClose(sessionW)
    // loading the shell
    sessionW.WriteTransaction(func(tx neo4j.Transaction) (interface{}, error) {
        loadShellTx(strId, fileName, shell, bounds) // updated to include fileName reference
        return nil, nil
    })
    // end loading the shell
    fmt.Fprint(w, "1. Shell thickness in Angstrom:%f\n", shellThickness)
    fmt.Fprint(w, "2. Maximum number of bounds per atom:%d\n", maxBounds)
    fmt.Fprint(w, "3. Size in Angstrom:%f\n", size)
    // no verbose fingerprint generation
    fmt.Fprint(w2, "%f\n", shellThickness)
    fmt.Fprint(w2, "%d\n", maxBounds)
    fmt.Fprint(w2, "%f\n", size)
    // not connected
    // read only transaction to make the queries:
    session := driver.NewSession(neo4j.SessionConfig{AccessMode: neo4j.AccessModeRead})
    defer unsafeClose(session)
    // doing the readonly queries
    session.ReadTransaction(func(tx neo4j.Transaction) (interface{}, error) {
        w.Flush()
        w2.Flush()
        sessionW.WriteTransaction(func(tx neo4j.Transaction) (interface{}, error) {
            deleteShell(tx, strId)
            return nil, nil
        })
    })
    return filepath.Join(strId, fileOut), filepath.Join(strId, fileOutNoVerbose), nil
}

```

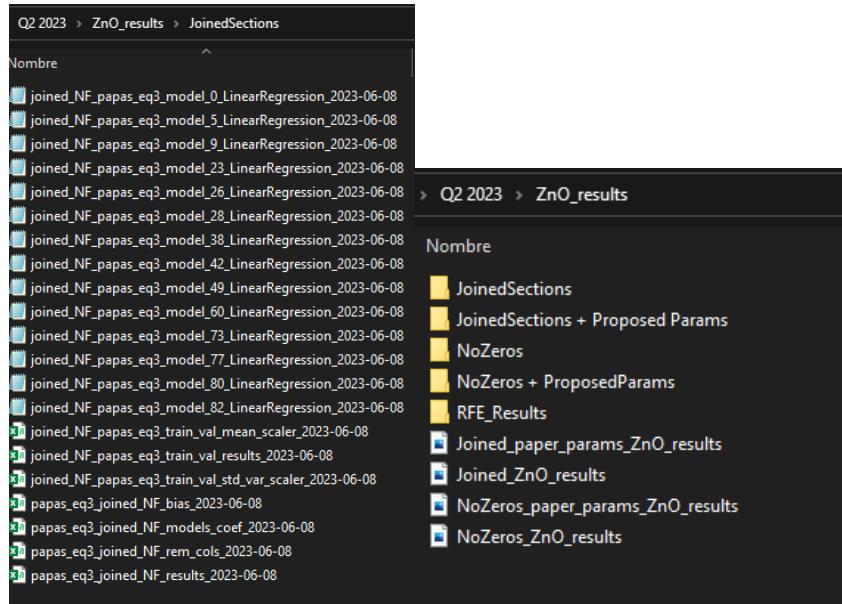
Some steps have been hidden to avoid oversharing contents of site backend – however, as seen in the snippets, GenVector function is in charge of performing the required operations, such as initializing files where the NanoFingerprint will be written to as well as computing what atoms are part of the shell, to later submit it and perform additional operations using Neo4j.

## Code model analysis

Code used for performing model analysis has been added to GitHub – it can be found at:  
<https://github.com/khuillca/NanoFingerprint-model-analysis>.

## Models Generated

A set of multiple models have been created for the experiments conducted in section 4. These models can be loaded on a python script to be analysed on further samples. These have been stored in 4 different directories, each of these have a set of results obtained from the models as well as the coefficients linked to each LinearRegression model generated.



All models and results have been shared with the director of the TFM, which will make them available to the public.

## Appendix B

### RFE Results

#### TiO2

	mean_fit_time	std_fit_time	mean_scoring_time	std_scoring_time	param_n_features_to_select	params	split0_test_r2	split1_test_r2	split2_test_r2	split3_test_r2	split4_test_r2	mean_test_r2	std_test_r2	rank_test_r2	split0_train_r2	split1_train_r2	split2_train_r2	split3_train_r2	split4_train_r2	mean_train_r2	std_train_r2
	0.000	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
0	0.0512 5518	5002 2	0.00300 288	3.37E-07 0038	1	{'n_features_to_select': 1}	0.055 48076	0.169 53452	0.143 36964	0.850 70966	0.041 07557	0.2520 3403	3688 5	0.303 35	0.0443 0053	0.0813 4835	0.0010 3552	0.1308 0274	0.0123 7219	0.0539 7187	0.047 4948
1	0.0500 6728	4.16E-05 252	0.00280 0004	0.0004 0038	2	{'n_features_to_select': 2}	1.009 55567	0.316 30052	0.018 1279	0.004 14716	0.346 96378	0.0664 6213	5366 8	0.493 34	0.4718 1351	0.4180 5398	0.5402 5596	0.4441 168	0.3949 922	0.4538 4649	0.2573 0057
2	0.0484 6296	5143 3	0.00260 234	0.0004 9037	3	{'n_features_to_select': 3}	0.306 33734	0.679 94614	0.242 43613	0.325 48503	0.195 32156	0.2529 3079	3665 6	0.296 17	0.6516 1896	0.5985 6408	0.7761 7703	0.6847 6495	0.6737 8899	0.6769 828	0.7980 0056
3	0.0470 439	2.34E-06 204	0.00280 0181	0.0004 0000	4	{'n_features_to_select': 4}	0.375 04934	0.678 49598	0.315 50648	0.236 38256	0.234 3456	0.2417 534	2217 2	0.6666 19	0.6263 2856	0.7832 2091	0.7122 9756	0.6956 3298	0.6968 0415	0.1518 1683	0.052 0052
4	0.0452 7907	3545 2	0.00280 404	0.0004 0115	5	{'n_features_to_select': 5}	0.375 04934	0.678 49598	0.315 50648	0.236 38256	0.150 17233	0.2249 1874	3607 2	0.324 27	0.6666 2856	0.6263 2091	0.7832 9756	0.7122 3298	0.7009 167	0.6978 7934	0.1703 0052
5	0.0448 8792	4298 5	0.00299 463	1.78E-05 005	6	{'n_features_to_select': 6}	0.375 04934	0.678 49598	0.315 50648	0.197 34089	0.150 17233	0.2171 1041	4606 7	0.324 33	0.6666 2856	0.6263 2091	0.7832 9756	0.7235 7721	0.7009 167	0.7001 4819	0.9855 0052
6	0.0430 9664	1147 5	0.00240 235	0.0004 9014	7	{'n_features_to_select': 7}	0.375 04934	0.678 49598	0.315 50648	0.197 34089	0.150 17233	0.2171 1041	4606 7	0.324 32	0.6666 2856	0.6263 2091	0.7832 9756	0.7235 7721	0.7009 167	0.7001 4819	0.9855 0051
7	0.0414 0525	5162 3	0.00251 398	0.0004 4833	8	{'n_features_to_select': 8}	0.375 04934	0.685 51919	0.315 50648	0.197 34089	0.150 17233	0.2185 1505	4640 6	0.326 29	0.6666 2856	0.6303 6488	0.7832 9756	0.7235 7721	0.7009 167	0.7009 5698	0.8716 0051
8	0.0420 3839	0533 8	0.00260 234	0.0004 9017	9	{'n_features_to_select': 9}	0.375 04934	0.685 51919	0.315 50648	0.197 34089	0.150 17233	0.2185 1505	4640 6	0.326 30	0.6666 2856	0.6303 6488	0.7832 9756	0.7235 7721	0.7009 167	0.7009 5698	0.8716 0051
9	0.0444 4041	4256 3	0.00280 299	0.0004 0026	10	{'n_features_to_select': 10}	0.375 04934	0.685 51919	0.315 50648	0.197 34089	0.150 17233	0.2185 1505	4640 6	0.326 28	0.6666 2856	0.6303 6488	0.7832 9756	0.7235 7721	0.7009 167	0.7009 5698	0.8716 0051
10	0.0376 3442	4905 1	0.00260 22	0.0004 9045	11	{'n_features_to_select': 11}	0.375 04934	0.685 51919	0.315 50648	0.197 34089	0.150 17233	0.2185 1505	4640 6	0.326 31	0.6666 2856	0.6303 6488	0.7832 9756	0.7235 7721	0.7009 167	0.7009 5698	0.8716 0060
11	0.0364 3308	4904 9	0.00260 248	0.0004 9049	12	{'n_features_to_select': 12}	0.375 04934	0.685 51919	0.247 31295	0.197 34089	0.150 17233	0.2321 5588	5376 9	0.304 26	0.6666 2856	0.6303 6488	0.8106 5227	0.7235 7721	0.7009 167	0.7064 2792	0.9208 007

1	0.0408	8712	0.00320	0.0014	{'n_features_to_select': 13}	0.375	0.685	0.247	0.197	0.150	0.2321	5588	-	0.304	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	0.060
2	3691	4	287	7105	{'n_features_to_select': 13}	04934	51919	31295	34089	17233	5376	9	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9208
1	0.0336	4904	0.00300	2.34E-0000	{'n_features_to_select': 14}	0.375	0.685	0.231	0.197	0.150	0.2352	7114	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	0.060
3	3061	7	269	07	{'n_features_to_select': 14}	04934	51919	84009	34089	17233	4833	1	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
1	0.0319	6361	0.00260	0.0004	{'n_features_to_select': 15}	0.375	0.685	0.231	0.197	0.150	0.2352	7114	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
4	9759	4	248	9	{'n_features_to_select': 15}	04934	51919	84009	34089	17233	4833	1	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
1	0.0310	6329	0.00200	3.57E-0000	{'n_features_to_select': 16}	0.375	0.685	0.231	0.197	0.150	0.2352	7114	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
5	2822	4	191	07	{'n_features_to_select': 16}	04934	51919	84009	34089	17233	4833	1	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
1	0.0298	4002	0.00280	0.0004	{'n_features_to_select': 17}	0.375	0.685	0.231	0.197	0.150	0.2352	7114	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
6	2698	3	285	0042	{'n_features_to_select': 17}	04934	51919	84009	34089	17233	4833	1	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
1	0.0281	6703	0.00260	0.0004	{'n_features_to_select': 18}	0.375	0.685	0.231	0.197	0.150	0.2352	7114	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
7	3549	2	229	9023	{'n_features_to_select': 18}	04934	51919	84009	34089	17233	4833	1	-	0.299	0.6666	0.6303	0.8106	0.7235	0.7009	0.7064	9367
1	0.0268	4005	0.00220	0.0004	{'n_features_to_select': 19}	0.375	0.685	0.231	0.238	0.150	0.2434	1230	-	0.299	0.6666	0.6303	0.8106	0.7615	0.7009	0.7140	8430
8	2452	7	189	005	{'n_features_to_select': 19}	04934	51919	84009	14418	17233	0899	8	-	0.310	0.6666	0.6303	0.8106	0.7615	0.7009	0.7140	8430
1	0.0255	4479	0.00252	0.0004	{'n_features_to_select': 20}	0.472	0.685	0.231	0.238	0.150	0.2629	0581	-	0.310	0.6983	0.6303	0.8106	0.7615	0.7009	0.7203	3510
9	3134	5	361	4963	{'n_features_to_select': 20}	62627	51919	84009	14418	17233	2438	3	-	0.310	0.6983	0.6303	0.8106	0.7615	0.7009	0.7203	3510
2	0.0242	4003	0.00280	0.0004	{'n_features_to_select': 21}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
0	2209	3	247	0047	{'n_features_to_select': 21}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0231	4471	0.00280	0.0004	{'n_features_to_select': 22}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
1	5946	1	271	0035	{'n_features_to_select': 22}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0214	4902	0.00280	0.0004	{'n_features_to_select': 23}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	1962	9	228	005	{'n_features_to_select': 23}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0206	4903	0.00240	0.0004	{'n_features_to_select': 24}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
3	1868	7	221	9035	{'n_features_to_select': 24}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0188	2372	0.00280	0.0004	{'n_features_to_select': 25}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
4	2348	3	261	0042	{'n_features_to_select': 25}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0174	4904	0.00280	0.0004	{'n_features_to_select': 26}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
5	1571	1	261	0031	{'n_features_to_select': 26}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0165	4705	0.00280	0.0003	{'n_features_to_select': 27}	0.472	0.714	0.231	0.238	0.150	0.2686	9431	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
6	9336	5	17	9986	{'n_features_to_select': 27}	62627	04319	84009	14418	17233	2918	7	-	0.310	0.6983	0.6585	0.8106	0.7615	0.7009	0.7260	6335
2	0.0150	2.28E-05	0.00297	4.84E-0000	{'n_features_to_select': 28}	0.422	0.714	0.231	0.238	0.150	0.2586	0885	-	0.310	0.6986	0.6585	0.8106	0.7615	0.7009	0.7260	6019
7	0392	-05	78	05	{'n_features_to_select': 28}	49463	04319	84009	14418	17233	0285	9	-	0.310	0.6986	0.6585	0.8106	0.7615	0.7009	0.7260	6019

ZnO

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_n_features_to_se lect	params	split0_test_r2	split1_test_r2	mean_test_r2	std_test_r2	rank_test_r2	split0_train_r2	split1_train_r2	mean_train_r2	std_train_r2
0	0.19766223	0.000857	0.0048902	0.00011349	1	{'n_features_to_select': 1}	2.66891892	28.4637097	15.5663143	12.89739	141	0.60483871	0.66216216	0.63350044	0.0286617
1	0.19502616	0.0001509	0.00500441	2.38E-07	2	{'n_features_to_select': 2}	7.27117586	9.34301241	8.30709414	1.035918	8	1	0.97838343	0.98919172	0.0108082
2	0.19424963	0.0017032	0.00491035	9.60E-05	3	{'n_features_to_select': 3}	9.76684136	-10.733022	10.2499317	0.483090	19	1	1	1	0
3	0.19076061	0.0004794	0.00500596	1.19E-07	4	{'n_features_to_select': 4}	-9.4800264	13.4409937	11.4605101	1.980483	134	1	1	1	0
4	0.18883145	0.0013417	0.00450444	0.00050044	5	{'n_features_to_select': 5}	3.15946718	17.2781005	10.2187838	7.059316	16	1	1	1	0
5	0.18786907	0.001302	0.00500512	2.38E-07	6	{'n_features_to_select': 6}	3.15844905	17.2893053	10.2238772	7.065428	17	1	1	1	0
6	0.187717015	0.0005013	0.00500453	3.58E-07	7	{'n_features_to_select': 7}	3.15664666	17.2984046	10.2275256	7.070878	18	1	1	1	0
7	0.18466902	0.0001732	0.0045042	0.00050044	8	{'n_features_to_select': 8}	3.15723848	17.4613791	10.3093088	7.152070	20	1	1	1	0
8	0.18342841	0.0010006	0.00450456	0.00050056	9	{'n_features_to_select': 9}	3.15779378	17.4651761	10.3114849	7.153691	21	1	1	1	0
9	0.18216562	0.0005012	0.00500464	2.38E-07	10	{'n_features_to_select': 10}	3.15932165	17.4649087	10.3121152	7.152793	22	1	1	1	0
10	0.18166602	0.0005012	0.00500441	0	11	{'n_features_to_select': 11}	3.15852632	17.4671431	10.3128347	7.154308	23	1	1	1	0
11	0.18016422	1.19E-07	0.00500441	0	12	{'n_features_to_select': 12}	3.15881715	17.4689934	10.3139053	7.155088	24	1	1	1	0

12	0.17825747	9.51E-05 0.0015015	0.00500441	0	13	{'n_features_to_select': 13} - 3.15914116 'n_features_to_select': - 17.4701391 14 14} - 10.3146401 'n_features_to_select': - 97 15 15} - 7.155498 'n_features_to_select': - 3.15903948 'n_features_to_select': - 17.4711107 16 16} - 10.3150751 'n_features_to_select': - 61 17 17} - 7.156035 'n_features_to_select': - 3.15893904 'n_features_to_select': - 17.4720043 18 18} - 10.3154716 'n_features_to_select': - 6 19 19} - 7.156532 'n_features_to_select': - 3.15765552 'n_features_to_select': - 17.4728727 20 20} - 10.3152641 'n_features_to_select': - 58 21 21} - 7.157608 'n_features_to_select': - 3.15811812 'n_features_to_select': - 17.4736911 22 22} - 10.3159046 'n_features_to_select': - 5 23 23} - 7.157786 'n_features_to_select': - 3.15739465 'n_features_to_select': - 17.4743644 24 24} - 10.3158795 'n_features_to_select': - 89 25 25} - 7.158017 'n_features_to_select': - 3.15677561 'n_features_to_select': - 17.4728111 26 26} - 10.3147934 'n_features_to_select': - 76 27 27} - 7.158073 'n_features_to_select': - 3.15704546 'n_features_to_select': - 17.4731931 28 28} - 10.3151193 'n_features_to_select': - 81 29 29} - 7.158484 'n_features_to_select': - 3.15697843 'n_features_to_select': - 17.4735598 30 30} - 10.3152691 'n_features_to_select': - 68 31 31} - 7.158290 'n_features_to_select': - 3.15691159 'n_features_to_select': - 17.4738827 32 32} - 10.3153972 'n_features_to_select': - 58 33 33} - 7.158485 'n_features_to_select': - 3.15714594 'n_features_to_select': - 17.4741941 34 34} - 10.3151567 'n_features_to_select': - 08 35 35} - 7.158524 'n_features_to_select': - 3.15738015 'n_features_to_select': - 17.4744668 36 36} - 10.3159235 'n_features_to_select': - 31 37 37} - 7.158543 'n_features_to_select': - 3.15745891 'n_features_to_select': - 17.4746332 38 38} - 10.316046 'n_features_to_select': - 13 39 39} - 7.158653 'n_features_to_select': - 3.15748681 'n_features_to_select': - 17.4747946 40 40} - 10.3161407 'n_features_to_select': - 9 41 41} - 7.158653 'n_features_to_select': - 4.44373391 'n_features_to_select': - 17.4750155 42 42} - 10.9593747 'n_features_to_select': - 78 43 43} - 7.158642 'n_features_to_select': - 3.15794394 'n_features_to_select': - 17.4752299 44 44} - 10.3165869 'n_features_to_select': - 99 45 45} - 7.158743 'n_features_to_select': - 3.15800895 'n_features_to_select': - 17.4754964 46 46} - 10.3167527 'n_features_to_select': - 71 47 47} - 7.158692 'n_features_to_select': - 3.15826389 'n_features_to_select': - 17.475649 48 48} - 10.3169565 'n_features_to_select': - 56 49 49} - 7.158675 'n_features_to_select': - 3.15844002 'n_features_to_select': - 17.4757907 50 50} - 10.3171154 'n_features_to_select': - 34 51 51} - 6.215553 'n_features_to_select': - 3.15853698 'n_features_to_select': - 15.5896447 52 52} - 9.37409083 'n_features_to_select': - 85 53 53} - 7.159067 'n_features_to_select': - 3.15846237 'n_features_to_select': - 17.4765973 54 54} - 10.3175298 'n_features_to_select': - 46 55 55} - 7.159074 'n_features_to_select': - 3.1585459 'n_features_to_select': - 17.4766954 56 56} - 10.3176207 'n_features_to_select': - 75 57 57} - 7.159075 'n_features_to_select': - 3.1586294 'n_features_to_select': - 17.4767804 58 58} - 10.3177049 'n_features_to_select': - 48 59 59} - 6.512568 'n_features_to_select': - 3.15867004 'n_features_to_select': - 16.1838077 60 60} - 9.67123885 'n_features_to_select': - 81 61 61} - 7.159136 'n_features_to_select': - 3.15871068 'n_features_to_select': - 17.4769844 62 62} - 10.3178476 'n_features_to_select': - 87 63 63} - 7.159179 'n_features_to_select': - 3.15870757 'n_features_to_select': - 17.4770675 64 64} - 10.3178875 'n_features_to_select': - 96 65 65} - 7.159275 'n_features_to_select': - 3.15869571 'n_features_to_select': - 17.4772466 66 66} - 10.3179712 'n_features_to_select': - 46 67 67} - 6.667238 'n_features_to_select': - 3.15874129 'n_features_to_select': - 16.4932189 68 68} - 9.82598008 'n_features_to_select': - 8 69 69} - 15 'n_features_to_select': - 1 'n_features_to_select': - 1 'n_features_to_select': - 1 'n_features_to_select': - 1 'n_features_to_select': - 0
----	------------	-----------------------	------------	---	----	---









## TiO<sub>2</sub>+ZnO

	mean_fit_time	std_fit_time	mean_scorer_time	std_scoring_time	param_n_features_to_select	params	split0_t est_r2	split1_t est_r2	split2_t est_r2	split3_t est_r2	split4_t est_r2	mean_t est_r2	std_tes t_r2	rank_t est_r2	split0_tr ain_r2	split1_tr ain_r2	split2_tr ain_r2	split3_tr ain_r2	split4_tr ain_r2	mean_t rain_r2	std_tra in_r2
0	0.24389 0.0572	0.0021 20067	0.005316 353	0.00040 6627	1	{'n_features_to_select': 1}	0.0021 94267	0.0638 68588	0.2057 2061	0.3624 34142	3.1651 7589	0.6511 65313	1.2709 06808	135	0.09214 9804	0.08868 6557	0.00422 9603	0.00205 504	0.18672 1209	0.07476 8443	0.0682 52787
1	0.24108 6054	0.0015 31413	0.005805 063	0.00040 0209	2	{'n_features_to_select': 2}	0.0032 5219	0.1505 94031	0.2153 79221	0.1613 19306	3.1217 39333	0.5827 6664	1.2761 0916	122	0.10025 0945	0.09686 0741	0.00490 2291	0.17488 5014	0.19484 913	0.11434 9624	0.0672 72664
2	0.23879 8475	0.0021 43931	0.005404 806	0.00049 0388	3	{'n_features_to_select': 3}	0.0032 5219	0.1174 28132	0.8611 41059	0.1248 91825	3.4402 65664	0.8611 23645	1.3340 70464	165	0.10025 0945	0.13196 121	0.03081 6936	0.19545 5945	0.21478 8456	0.13465 4698	0.0664 7009
3	0.23697 052	0.0013 80391	0.005004 692	3.81E-07	4	{'n_features_to_select': 4}	0.0252 75114	0.1485 8371	0.8611 41059	0.1255 95399	5.0887 11067	1.1803 1774	1.9854 11018	166	0.17205 5916	0.13272 1231	0.03081 6936	0.19545 7474	0.30556 8196	0.16732 395	0.0891 70479
4	0.23521 4329	0.0018 99039	0.005204 535	0.00040 0019	5	{'n_features_to_select': 5}	0.4435 01962	0.1485 8371	2.3843 75948	0.1264 06496	5.0887 11067	1.4014 81568	2.0973 42403	168	0.49117 9623	0.13272 1231	0.27907 7486	0.19548 231	0.30556 8196	0.28080 5769	0.1217 23002
5	0.23418 3073	0.0014 26241	0.005204 391	0.00040 0448	6	{'n_features_to_select': 6}	0.4458 16329	0.1485 8371	2.1395 13878	0.1264 06496	2.2078 51111	0.7758 74289	1.1557 70697	163	0.49382 8241	0.13272 1231	0.29349 9863	0.19548 231	0.53461 8149	0.33002 9959	0.1594 0794
6	0.24194 6554	0.0150 87192	0.005404 902	0.00049 0408	7	{'n_features_to_select': 7}	0.4458 16329	0.1485 8371	2.1395 13878	0.1264 06496	1.2755 61018	0.5894 16271	0.9699 82037	124	0.49382 8241	0.13272 1231	0.29349 9863	0.19548 231	0.70154 3767	0.36341 5082	0.2086 37392
7	0.23091 507	0.0015 53096	0.005102 015	0.00049 0278	8	{'n_features_to_select': 8}	0.8176 38667	0.1485 8371	2.1395 13878	0.1264 06496	1.2755 61018	0.5150 51803	1.0568 60677	11	0.57254 7524	0.13272 1231	0.29349 9863	0.19548 231	0.70154 3767	0.37915 8939	0.2205 16325
8	0.23039 2122	0.0014 87383	0.005600 405	0.00048 6385	9	{'n_features_to_select': 9}	0.8176 38667	0.1485 8371	2.1395 13878	0.1264 06496	1.9561 43271	0.6511 68254	1.1824 13157	136	0.57254 7524	0.13272 1231	0.29349 9863	0.19548 231	0.81904 5	0.40265 9186	0.2568 82739
9	0.22960 4912	0.0017 5004	0.005398 846	0.00048 3647	10	{'n_features_to_select': 10}	0.8016 31883	0.1485 8371	2.1395 13878	0.1264 06496	1.9561 43271	0.6543 6961	1.1784 47079	139	0.66554 2308	0.13272 1231	0.29349 9863	0.19548 231	0.81904 5	0.42125 8142	0.2714 62554
0	0.22840 8813	0.0021 55364	0.005197 811	0.00039 9025	11	{'n_features_to_select': 11}	0.8016 31883	0.1485 8371	2.1395 13878	0.1264 06496	1.9561 43271	0.6543 6961	1.1784 47079	141	0.66554 2308	0.13272 1231	0.29349 9863	0.19548 231	0.81904 5	0.42125 8142	0.2714 62554
1	0.22780 1418	0.0017 40942	0.005186 224	0.00041 2283	12	{'n_features_to_select': 12}	0.8016 31883	0.1485 8371	2.1395 13878	0.1264 06496	1.9561 43271	0.6543 6961	1.1784 47079	140	0.66554 2308	0.13272 1231	0.29349 9863	0.19548 231	0.81904 5	0.42125 8142	0.2714 62554
2	0.22440 877	0.0018 63506	0.005200 672	0.00040 2366	13	{'n_features_to_select': 13}	0.8016 31883	0.3915 10774	2.1395 13878	0.1264 06496	1.9561 43271	0.6057 84198	1.2149 9108	128	0.66554 2308	0.63491 1647	0.29349 9863	0.19548 231	0.81904 5	0.52169 6226	0.2368 1842
3	0.22362 0987	0.0018 867	0.005404 854	0.00049 0155	14	{'n_features_to_select': 14}	0.8016 31883	0.3915 10774	2.1395 13878	0.1264 06496	1.9561 43271	0.6057 84198	1.2149 9108	127	0.66554 2308	0.63491 1647	0.29349 9863	0.19548 231	0.81904 5	0.52169 6226	0.2368 1842
4	0.22234 025	0.0018 6502	0.005805 54	0.00039 9971	15	{'n_features_to_select': 15}	0.8016 31883	0.3915 10774	2.1395 13878	0.1264 06496	1.9561 43271	0.6057 84198	1.2149 9108	126	0.66554 2308	0.63491 1647	0.29349 9863	0.19548 231	0.81904 5	0.52169 6226	0.2368 1842
5	0.22166 9674	0.0026 89131	0.005404 758	0.00049 0524	16	{'n_features_to_select': 16}	0.8016 31883	0.3930 18825	2.1395 13878	0.1264 06496	1.9561 43271	0.6054 82587	1.2152 38774	125	0.66554 2308	0.63493 3769	0.29349 9863	0.19548 231	0.81904 5	0.52170 065	0.2368 20535
6	0.21926 5842	0.0015 95432	0.005605 984	0.00049 0214	17	{'n_features_to_select': 17}	0.8016 31883	0.5620 81732	2.1395 13878	0.1264 06496	1.9561 43271	0.5716 70006	1.2445 49019	80	0.66554 2308	0.75871 9199	0.29349 9863	0.19548 231	0.81904 5	0.54645 7736	0.2532 63229



3	0.19290	0.0013	0.005004	3.57E-07	6	9908	12842	549		37	{'n_features_to_select': 37}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
3	0.19077	0.0014	0.004904	0.00020	7	3773	9788	318	125	38	{'n_features_to_select': 38}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
3	0.18961	0.0007	0.005204	0.00040	8	5059	98382	105	0831	39	{'n_features_to_select': 39}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
3	0.18784	0.0014	0.005368	0.00051	9	3704	43359	71	9135	40	{'n_features_to_select': 40}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18695	0.0011	0.005402	0.00048	0	4021	3695	47	7638	41	{'n_features_to_select': 41}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18575	0.0012	0.005204	0.00040	1	4824	18351	582	0591	42	{'n_features_to_select': 42}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18406	0.0014	0.005204	0.00040	2	7965	97198	725	0162	43	{'n_features_to_select': 43}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18330	0.0015	0.005404	0.00049	3	5979	35293	997	0427	44	{'n_features_to_select': 44}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18155	0.0008	0.005204	0.00040	4	8371	66641	773	0496	45	{'n_features_to_select': 45}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.18027	0.0012	0.005204	0.00040	5	1006	01815	868	009	46	{'n_features_to_select': 46}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.17816	0.0010	0.005404	0.00049	6	2527	96889	806	0291	47	{'n_features_to_select': 47}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.17736	0.0011	0.005004	6.14E-07	7	1488	67897	549		48	{'n_features_to_select': 48}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.17620	0.0012	0.005004	4.62E-07	8	5683	6942	358	07	49	{'n_features_to_select': 49}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
4	0.17474	0.0008	0.005205	0.00040	9	9994	15453	202	1235	50	{'n_features_to_select': 50}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
5	0.17363	0.0018	0.005204	0.00040	0	863	01251	964	1354	51	{'n_features_to_select': 51}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
5	0.17264	0.0014	0.005403	0.00048	1	6618	18058	423	8714	52	{'n_features_to_select': 52}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
5	0.17141	0.0013	0.005201	0.00039	2	9621	88339	435	8365	53	{'n_features_to_select': 53}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
5	0.16956	0.0016	0.005004	4.10E-07	3	8014	06471	597		54	{'n_features_to_select': 54}	0.7998	0.5132	2.1395	0.1264	1.9561	0.5818	1.2353	-	0.68138	0.75987	0.29349	0.19548	0.81904	0.54985	0.2550	7324	17257
5	0.16808	0.0013	0.005004	5.35E-07	4	9104	69845	501		55	{'n_features_to_select': 55}	0.7998	0.5132	2.1395	0.1264	1.7816	0.5469	1.1979	-	0.68138	0.75987	0.29349	0.19548	0.82227	0.55050	0.2557	3415	01605





9	0.11370	0.0008	0.005004	7.89E-07	3346	00658	358	07	94	{'n_features_to_select': 94}	0.7998	0.5132	2.1395	0.1264	1.7816	0.5469	1.1979	-	31	0.68138	0.75987	0.29349	0.19548	0.82227	0.55050	0.2557
9	0.11210	0.0006	0.005611	0.00048	2318	32938	849	2368	95	{'n_features_to_select': 95}	0.7998	0.5132	2.1395	0.1264	1.7816	0.5469	1.1979	-	23	0.68138	0.75987	0.29349	0.19548	0.82227	0.55050	0.2557
9	0.11063	0.0006	0.005004	4.77E-07	1752	60218	883	07	96	{'n_features_to_select': 96}	0.7998	0.5132	2.1395	0.1194	1.7816	0.5455	1.1984	-	13	0.68138	0.75987	0.29349	0.19536	0.82227	0.55048	0.2557
9	0.10978	0.0011	0.005010	4.91E-06	7846	95224	462	06	97	{'n_features_to_select': 97}	0.7998	0.5132	2.9146	0.1266	1.7816	0.7019	1.4230	-	160	0.68138	0.75987	0.26815	0.19548	0.82227	0.54543	0.2609
9	0.10812	0.0010	0.005205	0.00040	149	77886	059	0949	98	{'n_features_to_select': 98}	0.6665	0.5132	2.1395	0.1491	1.7816	0.5781	1.1670	-	88	0.73242	0.75987	0.29349	0.19083	0.82227	0.55978	0.2629
9	0.10649	0.0010	0.005004	6.47E-07	6525	20273	358	07	99	{'n_features_to_select': 99}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	71	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
9	0.10589	0.0007	0.005186	0.00041	5567	2768	272	1177	100	{'n_features_to_select': 100}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	73	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.10392	0.0007	0.005205	0.00040	7231	60396	107	021	101	{'n_features_to_select': 101}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	75	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.10317	0.0013	0.005004	2.43E-07	4353	60942	025	07	102	{'n_features_to_select': 102}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	65	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.10093	0.0007	0.005004	4.91E-07	4029	64899	263	07	103	{'n_features_to_select': 103}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	76	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.09889	0.0004	0.005204	0.00040	0162	00353	725	052	104	{'n_features_to_select': 104}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	62	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.09854	0.0005	0.005004	3.50E-07	9414	71351	597	07	105	{'n_features_to_select': 105}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	67	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.09646	0.0005	0.005204	0.00040	0485	36721	582	0472	106	{'n_features_to_select': 106}	0.6913	0.5132	2.1395	0.1240	1.7816	0.5681	1.1742	-	58	0.73310	0.75987	0.29349	0.19546	0.82227	0.56084	0.2617
0	0.09725	0.0011	0.005351	0.00038	0175	43182	019	2707	107	{'n_features_to_select': 107}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	66	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
0	0.09427	0.0007	0.005204	0.00040	3329	5307	439	0543	108	{'n_features_to_select': 108}	0.6913	0.5132	2.1395	0.1586	1.7816	0.5750	1.1717	-	87	0.73310	0.75987	0.29349	0.18034	0.82227	0.55782	0.2660
0	0.11083	0.0204	0.006000	0.00110	9462	40907	376	1459	109	{'n_features_to_select': 109}	0.6913	0.5132	2.1395	0.1154	1.7816	0.5664	1.1749	-	57	0.73310	0.75987	0.29349	0.19527	0.82227	0.56080	0.2618
0	0.09393	0.0023	0.005604	0.00049	9352	28522	982	0466	110	{'n_features_to_select': 110}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	74	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
1	0.09126	0.0015	0.005404	0.00049	9922	599	663	0213	111	{'n_features_to_select': 111}	0.6913	0.5132	2.1395	0.1264	1.7816	0.5685	1.1740	-	69	0.73310	0.75987	0.29349	0.19548	0.82227	0.56084	0.2617
1	0.08903	0.0008	0.005112	0.00021	2364	53024	743	596	112	{'n_features_to_select': 112}	0.6913	0.5132	2.1395	0.0956	1.7816	0.5624	1.1764	-	56	0.73310	0.75987	0.29349	0.18487	0.82227	0.55872	0.2647



1	3	0.05979	0.0004	0.005004	2.43E-07	311	132	{'n_features_to_select': 132}	0.6913	0.5132	2.1395	0.4712	1.7816	0.6375	1.1560	-	-	0.73310	0.75987	0.29349	0.23021	0.82227	0.47570	0.3996
1	1	2328	24689	-	-	0085	53479	-	13878	93243	10848	72881	81505	-	-	130	9286	1147	9863	1421	5454	8866	58642	
1	3	0.05824	0.0003	0.005387	0.00046	068	133	{'n_features_to_select': 133}	0.6913	0.5132	1.7378	0.1344	1.7816	0.4898	1.0727	-	-	0.73310	0.75987	0.27362	0.19052	0.82227	0.55588	0.2672
2	9187	81325	-	-	-	0085	53479	-	48398	99868	10848	8091	32179	-	-	10	9286	1147	1152	8317	5454	1071	60568	
1	3	0.05685	0.0004	0.005307	0.00040	1591	134	{'n_features_to_select': 134}	0.6913	0.5132	2.1395	0.1264	2.4520	0.7026	1.3325	-	-	0.73310	0.75987	0.29349	0.19548	0.82091	0.56057	0.2614
3	1912	00448	-	-	-	0085	53479	-	13878	06496	89185	91199	6962	-	-	161	9286	1147	9863	231	9898	6501	91224	
1	3	0.05565	0.0004	0.005004	4.77E-07	644	135	{'n_features_to_select': 135}	0.6913	0.5132	2.1477	0.1979	2.2159	0.6714	1.2688	-	-	0.73310	0.75987	0.29344	0.18212	0.82413	0.55853	0.2659
4	0473	90271	-	-	-	0085	53479	-	19746	90273	99484	31188	41891	-	-	153	9286	1147	6912	4266	8198	7962	00664	
1	3	0.05408	0.0005	0.005004	5.76E-07	501	136	{'n_features_to_select': 136}	0.6913	0.5132	2.1395	0.0247	2.2159	0.6252	1.2865	-	-	0.73310	0.75987	0.29349	0.04376	0.82413	0.51337	0.3359
5	3109	80828	-	-	-	0085	53479	-	13878	7272	99484	37416	08598	-	-	129	9286	1147	9863	2496	8198	1199	65635	
1	3	0.05268	0.0004	0.005203	0.00040	0569	137	{'n_features_to_select': 137}	0.6913	0.5132	2.1395	0.1138	2.2159	0.6529	1.2736	-	-	0.73310	0.75987	0.29349	0.19051	0.82413	0.56022	0.2635
6	2352	52294	-	-	-	0085	53479	-	13878	76948	99484	67349	28947	-	-	137	9286	1147	9863	444	8198	6587	24547	
1	3	0.05124	0.0004	0.005004	4.42E-07	549	138	{'n_features_to_select': 138}	0.6913	0.5132	2.1395	0.1264	2.2159	0.6554	1.2725	-	-	0.73310	0.75987	0.29349	0.19548	0.82413	0.56122	0.2621
7	6643	00543	-	-	-	0085	53479	-	13878	06496	99484	73259	77695	-	-	146	9286	1147	9863	231	8198	0161	34436	
1	3	0.05024	0.0004	0.005004	5.35E-07	549	139	{'n_features_to_select': 139}	0.6913	0.5132	2.1395	0.1264	2.2159	0.6554	1.2725	-	-	0.73310	0.75987	0.29349	0.19548	0.82413	0.56122	0.2621
8	5953	00448	-	-	-	0085	53479	-	13878	06496	99484	73259	77695	-	-	143	9286	1147	9863	231	8198	0161	34436	
1	3	0.04816	0.0005	0.005404	0.00049	0447	140	{'n_features_to_select': 140}	0.6913	0.5132	2.1395	0.2376	2.2159	0.6777	1.2640	-	-	0.73310	0.75987	0.29349	0.15948	0.82413	0.55402	0.2723
9	0458	18277	-	-	-	0085	53479	-	13878	41479	99484	20255	78087	-	-	156	9286	1147	9863	4444	8198	0588	74946	
1	4	0.04664	0.0004	0.005404	0.00049	615	141	{'n_features_to_select': 141}	0.6913	0.5132	2.1395	0.0509	2.2159	0.6403	1.2791	-	-	0.73310	0.75987	0.29349	0.17920	0.82413	0.55796	0.2667
0	2733	90914	-	-	-	0085	53479	-	13878	25825	99484	77125	94781	-	-	131	9286	1147	9863	9783	8198	5655	15994	
1	4	0.04544	0.0004	0.005204	0.00040	582	142	{'n_features_to_select': 142}	0.6913	0.5132	2.1395	0.0703	2.2159	0.6442	1.2774	-	-	0.73310	0.75987	0.29349	0.18947	0.82413	0.56001	0.2638
1	1675	9033	-	-	-	0085	53479	-	13878	29806	99484	57921	28843	-	-	133	9286	1147	9863	8839	8198	9467	15291	
1	4	0.04444	0.0004	0.005017	2.51E-05	614	143	{'n_features_to_select': 143}	0.6913	0.5132	2.1395	0.1264	2.2159	0.6554	1.2725	-	-	0.73310	0.75987	0.29349	0.19548	0.82413	0.56122	0.2621
2	0746	90271	-	-	-	0085	53479	-	13878	06496	99484	73259	77695	-	-	142	9286	1147	9863	231	8198	0161	34436	
1	4	0.04283	0.0004	0.005204	0.00040	725	144	{'n_features_to_select': 144}	0.6913	0.5132	2.1395	0.0596	2.2159	0.6421	1.2783	-	-	0.73310	0.75987	0.29349	0.15298	0.82413	0.55271	0.2742
3	9193	00496	-	-	-	0085	53479	-	13878	10639	99484	14087	9886	-	-	132	9286	1147	9863	0787	8198	9856	64926	
1	4	0.04137	0.0004	0.005404	0.00048	997	145	{'n_features_to_select': 145}	0.6913	0.5132	2.1395	0.2206	2.2159	0.6743	1.2652	-	-	0.73310	0.75987	0.29349	0.16734	0.82413	0.55559	0.2701
4	826	24385	-	-	-	0085	53479	-	13878	91869	99484	30333	75876	-	-	155	9286	1147	9863	5558	8198	281	60277	
1	4	0.03995	0.0001	0.005204	0.00040	9552	146	{'n_features_to_select': 146}	0.6913	0.5132	2.1395	0.1264	2.2159	0.6554	1.2725	-	-	0.73310	0.75987	0.29349	0.19548	0.82413	0.56122	0.2621
5	9764	57901	-	-	-	0085	53479	-	13878	06496	99484	73259	77695	-	-	149	9286	1147	9863	231	8198	0161	34436	
1	4	0.03843	0.0004	0.005004	3.81E-07	692	147	{'n_features_to_select': 147}	0.6913	0.5132	2.1395	0.1516	2.2159	0.6605	1.2705	-	-	0.73310	0.75987	0.29349	0.19085	0.82413	0.56029	0.2634
6	5078	90583	-	-	-	0085	53479	-	13878	41829	99484	20325	17772	-	-	151	9286	1147	9863	5064	8198	4712	28989	
1	4	0.03692	0.0003	0.005238	0.00039	676	148	{'n_features_to_select': 148}	0.6913	0.5132	2.1395	0.1636	2.2159	0.6629	1.2695	-	-	0.73310	0.75987	0.29349	0.18092	0.82413	0.55830	0.2662
7	3218	92779	-	-	-	0085	53479	-	13878	19181	99484	15796	66996	-	-	152	9286	1147	9863	5488	8198	8796	29148	
1	4	0.03583	0.0004	0.005204	0.00040	535	149	{'n_features_to_select': 149}	0.6913	0.5132	2.1395	0.1264	2.2159	0.6554	1.2725	-	-	0.73310	0.75987	0.29349	0.19548	0.82413	0.56122	0.2621
8	2977	0083	-	-	-	0085	53479	-	13878	06496	99484	73259	77695	-	-	148	9286	1147	9863	231	8198	0161	34436	
1	4	0.03403	1.78E-07	0.005004	2.78E-07	454	150	{'n_features_to_select': 150}	0.6913	0.5132	2.1395	0.2502	2.2159	0.6802	1.2632	-	-	0.73310	0.75987	0.29349	0.16773	0.82413	0.55567	0.2699
9	0962	-	-	-	-	0085	53479	-	13878	28711	99484	37702	11388	-	-	157	9286	1147	9863	8232	8198	1345	93414	

