

Software Requirements Specification for OCRacle: Latin Alphabet Character Recognition

Phillip Tran

February 6, 2025

Contents

1	Reference Material	iv
1.1	Table of Units	iv
1.2	Table of Symbols	iv
1.3	Abbreviations and Acronyms	v
2	Introduction	1
2.1	Purpose of Document	1
2.2	Scope of Requirements	1
2.3	Characteristics of Intended Reader	1
2.4	Organization of Document	1
3	General System Description	1
3.1	System Context	2
3.2	User Characteristics	3
3.3	System Constraints	3
4	Specific System Description	3
4.1	Problem Description	4
4.1.1	Terminology and Definitions	4
4.1.2	Physical System Description	4
4.1.3	Goal Statements	5
4.2	Solution Characteristics Specification	5
4.2.1	Assumptions	7
4.2.2	Theoretical Models	7
4.2.3	General Definitions	8
4.2.4	Data Definitions	9
4.2.5	Instance Models	10
4.2.6	Input Data Constraints	11
4.2.7	Properties of a Correct Solution	12
5	Requirements	12
5.1	Functional Requirements	13
5.2	Nonfunctional Requirements	13
5.3	Rationale	13
6	Likely Changes	14
7	Unlikely Changes	14
8	Traceability Matrices and Graphs	14

Revision History

Date	Version	Notes
January 27, 2025	1.0	Initial document creation

1 Reference Material

This section records information for easy reference.

1.1 Table of Units

This document does not use SI (Système International d'Unités) as the unit system because the problem domain is not a physical system.

symbol	unit	SI
px	pixel	N/A

1.2 Table of Symbols

The table that follows summarizes the symbols used in this document along with their units. The symbols are listed in alphabetical order.

symbol	unit	description
$I_{n \times m}$	$pixel^2$	Matrix representing an input image of n rows and m columns provided by the end user
$P_{1 \times 26}$	N/A	Probability vector representing the output of the model, with 26 columns representing the probability of each character
P_{pred}	N/A	The character predicted by the model.
$T_{28 \times 28}$	$pixel^2$	Matrix representing a 28x28 pixel image from the training data provided by the technical user

1.3 Abbreviations and Acronyms

symbol	description
A	Assumption
DD	Data Definition
GD	General Definition
GS	Goal Statement
IM	Instance Model
LC	Likely Change
PS	Physical System Description
R	Requirement
SRS	Software Requirements Specification
OCRacle	Name of the Project
TM	Theoretical Model

2 Introduction

Researchers analyzing physical print documents such as newspapers, books, and letters often need a means of digitizing the text in these documents. This enables them to search and analyze the text data more efficiently. Especially in the case of historical documents, digitizing the text can help preserve the information contained in these documents.

Optical Character Recognition (OCR) is a technology that allows for the extraction of text information from scanned documents, images, and other optical formats where text may be present. This digitalization process enables researchers to use computer programs to find trends and patterns in the digitized text.

The following section describes the purpose of the document, the scope of requirements, characteristics of the intended reader, and the organizational roadmap of the document.

2.1 Purpose of Document

This document serves as a software requirements specification for the OCRacle project. This includes the general system description, problem description, solution characteristics, and requirements. This document will be used a reference for building the solution.

2.2 Scope of Requirements

This program is intended for classifying a single handwritten uppercase Latin alphabet character in an image. Handwritten cursive characters are not an expected input of this program. The handwritten character are assumed to be in an upright position when input into the system.

2.3 Characteristics of Intended Reader

The intended reader of this document should have an undergraduate level understanding of data science, machine learning, and linear algebra.

2.4 Organization of Document

After the brief introduction, the document details the general system description, problem description, solutions characteristic specification, and requirements.

3 General System Description

This section provides general information about the system. It identifies the interfaces between the system and its environment, describes the user characteristics and lists the system constraints.

3.1 System Context

We can take two perspectives on the system context. The first is the technical user, who interfaces with the OCRacle program to produce a model that can predict the character in an image. The second is the end user, who uses the OCRacle model to predict the character in an input image.



Figure 1: System Context of End User

- Technical User Responsibilities:
 - Provide an image dataset to train the model following A1, A2, A3, A4, and A5.
 - Run the program on a compatible platform.
 - Interface with the program.
- OCRacle Responsibilities:
 - Detect incompatible input.
 - Train the model on the image dataset provided by the user.
 - Output a trained model that can predict the character in an image.
- End User Responsibilities:
 - Provide an image following A5.
 - Run the program on a compatible platform.

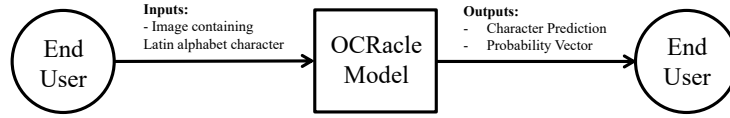


Figure 2: System Context of Technical User

- Interface with the program.
- OCRacle Responsibilities:
 - Detect incompatible input.
 - Provide an interface to view the predicted character.
 - Provide an interface to view the probability vector.

3.2 User Characteristics

Both the technical user and the end user of OCRacle should have an a basic understanding of the command line in order to setup the program. The provided user manual should be approachable for users of this skill level and allow them to train and use the model effectively, even if they are not familiar with the technical details of the program.

3.3 System Constraints

The system will operate solely on greyscale images of Latin uppercase characters and will not be able to predict characters in images containing multiple characters.

4 Specific System Description

This section first presents the problem description, which gives a high-level view of the problem to be solved. This is followed by the solution characteristics specification, which

presents the assumptions, theories, definitions and finally the instance models.

4.1 Problem Description

OCRacle is intended to solve the problem of extracting text information from a scanned document, image, and other optical formats where text may be present, such that this textual data can be used for further analysis.

4.1.1 Terminology and Definitions

This subsection provides a list of terms that are used in the subsequent sections and their meaning, with the purpose of reducing ambiguity and making it easier to correctly understand the requirements:

- **OCR:** Optical Character Recognition, the process of extracting text
- **EMNIST:** Extended MNIST, a dataset of handwritten characters
- **Latin Alphabet:** The alphabet used in the English language
- **Character:** A single letter in the Latin alphabet
- **Image:** A 2D array of pixel values
- **Pixel:** The smallest unit of a digital image
- **Probability Vector:** A vector representing the likelihood of each character
- **Model:** A trained machine learning model
- **Prediction:** The output of the model
- **Preprocessing:** The process of preparing the image for input into the model
- **Label:** The correct character associated with an image
- **Training:** The process of teaching the model to predict characters
- **Training Data:** The dataset used to train the model

4.1.2 Physical System Description

The physical system of OCRacle, as shown in Figure 3, includes the following elements:

PS1: For training the model, the system requires a dataset of images of Latin alphabet characters, alongside a corresponding label for each image.

PS2: For using the model as an end user, the system requires an image of a single Latin alphabet character.

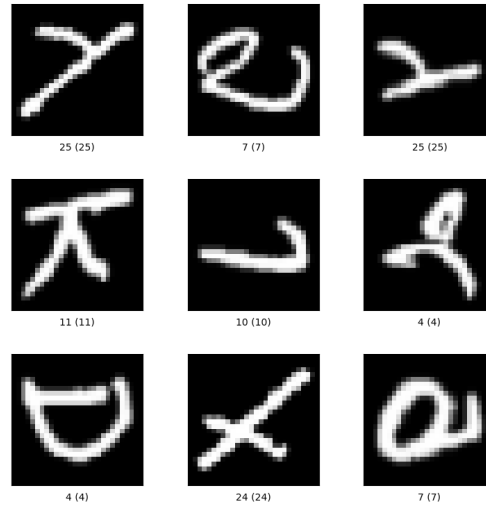


Figure 3: A sample of the EMNIST dataset, representing the Latin alphabet

4.1.3 Goal Statements

Given the training data provided by the technical user, the goal statements are:

GS1: Produce a trained model that can predict the character in an image.

Given an image containing a single Latin alphabet character from the end user, the goal statements are:

GS2: Given an image containing a single Latin alphabet character, predict the character.

GS3: Given an image containing a single Latin alphabet character, produce a probability vector representing the likelihood of each character.

4.2 Solution Characteristics Specification

[This section specifies the information in the solution domain of the system to be developed. This section is intended to express what is required in such a way that analysts and stakeholders get a clear picture, and the latter will accept it. The purpose of this section is to reduce the problem into one expressed in mathematical terms. Mathematical expertise is used to extract the essentials from the underlying physical description of the problem, and to collect and substantiate all physical data pertinent to the problem. —TPLT]

[This section presents the solution characteristics by successively refining models. It starts with the abstract/general Theoretical Models (TMs) and refines them to the concrete/specific

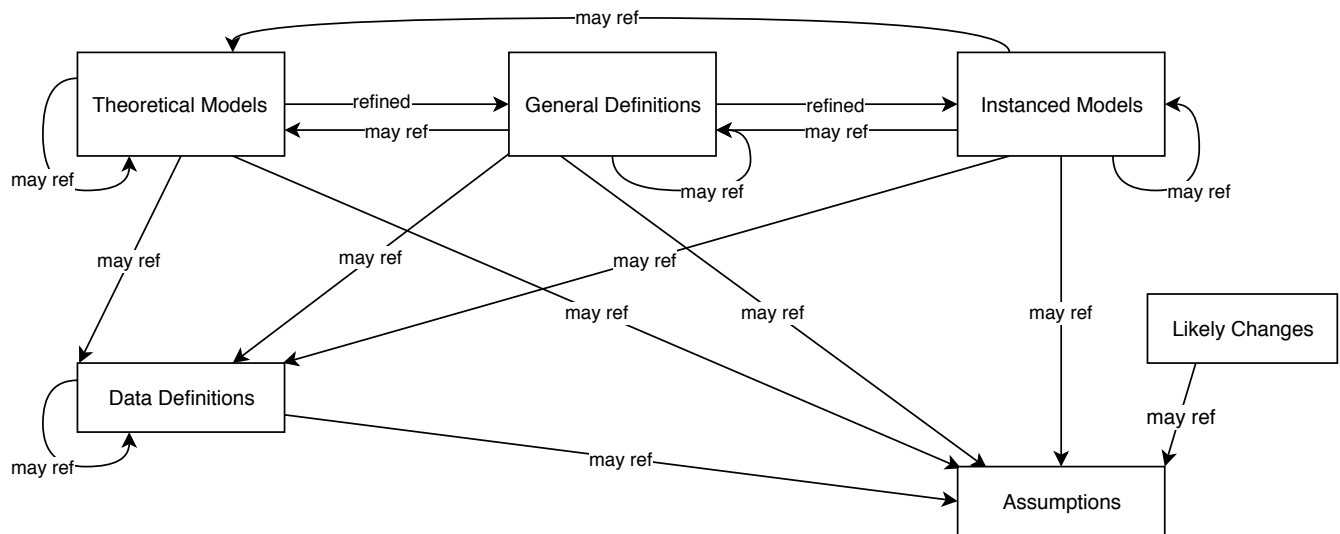
Instance Models (IMs). If necessary there are intermediate refinements to General Definitions (GDs). All of these refinements can potentially use Assumptions (A) and Data Definitions (DD). TMs are refined to create new models, that are called GMs or IMs. DDs are not refined; they are just used. GDs and IMs are derived, or refined, from other models. DDs are not derived; they are just given. TMs are also just given, but they are refined, not used. If a potential DD includes a derivation, then that means it is refining other models, which would make it a GD or an IM. —TPLT]

[The above makes a distinction between “refined” and “used.” A model is refined to another model if it is changed by the refinement. When we change a general 3D equation to a 2D equation, we are making a refinement, by applying the assumption that the third dimension does not matter. If we use a definition, like the definition of density, we aren’t refining, or changing that definition, we are just using it. —TPLT]

[The same information can be a TM in one problem and a DD in another. It is about how the information is used. In one problem the definition of acceleration can be a TM, in another it would be a DD. —TPLT]

[There is repetition between the information given in the different chunks (TM, GDs etc) with other information in the document. For instance, the meaning of the symbols, the units etc are repeated. This is so that the chunks can stand on their own when being read by a reviewer/user. It also facilitates reuse of the models in a different context. —TPLT]

[The relationships between the parts of the document are show in the following figure. In this diagram “may ref” has the same role as “uses” above. The figure adds “Likely Changes,” which are able to reference (use) Assumptions. —TPLT]



The instance models that govern OCRacle are presented in Subsection 4.2.5. The information to understand the meaning of the instance models and their derivation is also presented, so that the instance models can be verified.

4.2.1 Assumptions

This section simplifies the original problem and helps in developing the theoretical model by filling in the missing information for the physical system. The numbers given in the square brackets refer to the theoretical model [TM], general definition [GD], data definition [DD], instance model [IM], or likely change [LC], in which the respective assumption is used.

- A1: The model will be trained on the EMNIST dataset.
- A2: The training dataset will only contain images of Latin uppercase characters.
- A3: The training dataset will have one character per image.
- A4: The training dataset contains correct labels for each character represented in the image.
- A5: All input images from the user undergo preprocessing such that they are in the same format as the training data.

4.2.2 Theoretical Models

[Theoretical models are sets of abstract mathematical equations or axioms for solving the problem described in Section “Physical System Description” (Section 4.1.2). Examples of theoretical models are physical laws, constitutive equations, relevant conversion factors, etc. —TPLT]

[Optionally the theory section could be divided into subsections to provide more structure and improve understandability and reusability. Potential subsections include the following: Context theories, background theories, helper theories, generic theories, problem specific theories, final theories and rationale theories. —TPLT]

This section focuses on the general equations and laws that OCRacle is based on. [Modify the examples below for your problem, and add additional models as appropriate. —TPLT]

RefName: TM:COE

Label: Conservation of thermal energy

Equation: $-\nabla \cdot \mathbf{q} + g = \rho C \frac{\partial T}{\partial t}$

Description: The above equation gives the conservation of energy for transient heat transfer in a material of specific heat capacity C ($\text{J kg}^{-1} \text{°C}^{-1}$) and density ρ (kg m^{-3}), where \mathbf{q} is the thermal flux vector (W m^{-2}), g is the volumetric heat generation (W m^{-3}), T is the temperature (°C), t is time (s), and ∇ is the gradient operator. For this equation to apply, other forms of energy, such as mechanical energy, are assumed to be negligible in the system (A??). In general, the material properties (ρ and C) depend on temperature.

Notes: None.

Source: http://www.efunda.com/formulae/heat_transfer/conduction/overview_cond.cfm

Ref. By: GD??

Preconditions for TM:COE: None

Derivation for TM:COE: Not Applicable

[“Ref. By” is used repeatedly with the different types of information. This stands for Referenced By. It means that the models, definitions and assumptions listed reference the current model, definition or assumption. This information is given for traceability. Ref. By provides a pointer in the opposite direction to what we commonly do. You still need to have a reference in the other direction pointing to the current model, definition or assumption. As an example, if TM1 is referenced by GD2, that means that GD2 will explicitly include a reference to TM1. —TPLT]

4.2.3 General Definitions

[General Definitions (GDs) are a refinement of one or more TMs, and/or of other GDs. The GDs are less abstract than the TMs. Generally the reduction in abstraction is possible through invoking (using/referencing) Assumptions. For instance, the TM could be Newton’s

Law of Cooling stated abstracting. The GD could take the general law and apply it to get a 1D equation. —TPLT]

This section collects the laws and equations that will be used in building the instance models.

[Some projects may not have any content for this section, but the section heading should be kept. —TPLT] [Modify the examples below for your problem, and add additional definitions as appropriate. —TPLT]

Number	GD1
Label	Newton's law of cooling
SI Units	W m^{-2}
Equation	$q(t) = h\Delta T(t)$
Description	<p>Newton's law of cooling describes convective cooling from a surface. The law is stated as: the rate of heat loss from a body is proportional to the difference in temperatures between the body and its surroundings.</p> <p>$q(t)$ is the thermal flux (W m^{-2}).</p> <p>h is the heat transfer coefficient, assumed independent of T (A??) ($\text{W m}^{-2} \text{ }^{\circ}\text{C}^{-1}$).</p> <p>$\Delta T(t) = T(t) - T_{\text{env}}(t)$ is the time-dependent thermal gradient between the environment and the object ($^{\circ}\text{C}$).</p>
Source	Citation here
Ref. By	DD1, DD??

Detailed derivation of simplified rate of change of temperature

[This may be necessary when the necessary information does not fit in the description field. —TPLT] [Derivations are important for justifying a given GD. You want it to be clear where the equation came from. —TPLT]

4.2.4 Data Definitions

[The Data Definitions are definitions of symbols and equations that are given for the problem. They are not derived; they are simply used by other models. For instance, if a problem depends on density, there may be a data definition for the equation defining density. The DDs are given information that you can use in your other modules. —TPLT]

[All Data Definitions should be used (referenced) by at least one other model. —TPLT]

This section collects and defines all the data needed to build the instance models. The dimension of each quantity is also given. [Modify the examples below for your problem, and add additional definitions as appropriate. —TPLT]

Number	DD1
Label	Heat flux out of coil
Symbol	q_C
SI Units	W m^{-2}
Equation	$q_C(t) = h_C(T_C - T_W(t))$, over area A_C
Description	T_C is the temperature of the coil ($^{\circ}\text{C}$). T_W is the temperature of the water ($^{\circ}\text{C}$). The heat flux out of the coil, q_C (W m^{-2}), is found by assuming that Newton’s Law of Cooling applies (A??). This law (GD1) is used on the surface of the coil, which has area A_C (m^2) and heat transfer coefficient h_C ($\text{W m}^{-2}^{\circ}\text{C}^{-1}$). This equation assumes that the temperature of the coil is constant over time (A??) and that it does not vary along the length of the coil (A??).
Sources	Citation here
Ref. By	IM1

Type Name	Name for Type
Type Def	mathematical definition of the type
Description	description here
Sources	Citation here, if the type is borrowed from another source

4.2.5 Instance Models

This section transforms the problem defined in Section 4.1 into one which is expressed in mathematical terms. It uses concrete symbols defined in Section 4.2.4 to replace the abstract symbols in the models identified in Sections 4.2.2 and 4.2.3.

The goals [reference your goals —TPLT] are solved by [reference your instance models —TPLT]. [other details, with cross-references where appropriate. —TPLT] [Modify the examples below for your problem, and add additional models as appropriate. —TPLT]

Number	IM1
Label	ADAM for Minimizing Loss Function
Input	$m_W, C_W, h_C, A_C, h_P, A_P, t_{\text{final}}, T_C, T_{\text{init}}, T_P(t)$ from IM?? The input is constrained so that $T_{\text{init}} \leq T_C$ (A??)
Output	$T_W(t)$, $0 \leq t \leq t_{\text{final}}$, such that $\frac{dT_W}{dt} = \frac{1}{\tau_W}[(T_C - T_W(t)) + \eta(T_P(t) - T_W(t))]$, $T_W(0) = T_P(0) = T_{\text{init}}$ (A??) and $T_P(t)$ from IM??
Description	T_W is the water temperature ($^{\circ}\text{C}$). T_P is the PCM temperature ($^{\circ}\text{C}$). T_C is the coil temperature ($^{\circ}\text{C}$). $\tau_W = \frac{m_W C_W}{h_C A_C}$ is a constant (s). $\eta = \frac{h_P A_P}{h_C A_C}$ is a constant (dimensionless). The above equation applies as long as the water is in liquid form, $0 < T_W < 100^{\circ}\text{C}$, where 0°C and 100°C are the melting and boiling points of water, respectively (A??, A??).
Sources	Citation here
Ref. By	IM??

Derivation of ...

[The derivation shows how the IM is derived from the TMs/GDs. In cases where the derivation cannot be described under the Description field, it will be necessary to include this subsection. —TPLT]

4.2.6 Input Data Constraints

Table 1 shows the data constraints on the input output variables. The column for physical constraints gives the physical limitations on the range of values that can be taken by the variable. The column for software constraints restricts the range of inputs to reasonable values. The software constraints will be helpful in the design stage for picking suitable algorithms. The constraints are conservative, to give the user of the model the flexibility to experiment with unusual situations. The column of typical values is intended to provide a feel for a common scenario. The uncertainty column provides an estimate of the confidence with which the physical quantities can be measured. This information would be part of the input if one were performing an uncertainty quantification exercise.

The specification parameters in Table 1 are listed in Table 2.

Table 1: Input Variables

Var	Physical Constraints	Software Constraints	Typical Value	Uncertainty
$I_{n \times m}$	N/A	$0 \leq I_{ij} \leq 255$	N/A	N/A
$T_{28 \times 28}^*$	N/A	$0 \leq I_{ij} \leq 255$	N/A	N/A

(*) The training images are restricted to a 28x28 pixel format because the EMNIST dataset is in this format.

Table 2: Specification Parameter Values

Var	Value
n_{min}	1
n_{max}	2048
m_{min}	1
m_{max}	2048

4.2.7 Properties of a Correct Solution

As reflected in Table 3, the correct solution from the model must exhibit the following properties:

- The resulting probability vector must sum to 1.
- The predicted character must be one of the uppercase Latin alphabet characters from A to Z.

5 Requirements

This section provides the functional requirements, the business tasks that the software is expected to complete, and the nonfunctional requirements, the qualities that the software is expected to exhibit.

Table 3: Output Variables

Var	Physical Constraints
$P_{1 \times 26}$	$0 \leq P_i \leq 1$
P_{pred}	$P_{pred} \in \{A, B, C, \dots, Z\}$

5.1 Functional Requirements

- R1: The program accepts an image from the user in JPEG or PNG format.
- R2: The input image is processed such that it can be used for classification by the program.
- R3: The OCRacle project produces a model that predicts the character in an input image.
- R4: The model produced by the program outputs a probability vector to verify the prediction in a human readable-format.
- R5: The program outputs the most likely character prediction in a human-readable format.

5.2 Nonfunctional Requirements

- NFR1: **Accuracy** The accuracy of the the software shall exceed the previous OAR project. Since the OAR project provides a confusion matrix, the accuracy of OCRacle shall be measured by the same metric. [Ceranica \(2024\)](#)
- NFR2: **Usability** A user with a basic understand of the command line should be able to setup and run the program exclusively using the included user manual.
- NFR3: **Maintainability** The code should be highly modular, such that each processing step can be easily understood and re-implemented without disrupting the other previous and following processing steps.
- NFR4: **Portability** The program will be compatible with Windows, MacOS, and Linux operation systems. Any modern computers capable of running the operating systems mentioned above should be able to run the program.

5.3 Rationale

[Provide a rationale for the decisions made in the documentation. Rationale should be provided for scope decisions, modelling decisions, assumptions and typical values. —TPLT]

6 Likely Changes

LC1: The program may be modified to work non-Latin alphabet characters instead. For instance, Chinese character recognition may become the focus of the program.

LC2: The program may be extended to recognize lowercase Latin alphabet characters

LC3: The program may be extended to recognize Latin alphabet punctuation characters.

7 Unlikely Changes

UC1: The program will not be expanded to recognize full words using Latin alphabet characters.

8 Traceability Matrices and Graphs

The purpose of the traceability matrices is to provide easy references on what has to be additionally modified if a certain component is changed. Every time a component is changed, the items in the column of that component that are marked with an “X” may have to be modified as well. Table 4 shows the dependencies of theoretical models, general definitions, data definitions, and instance models with each other. Table 5 shows the dependencies of instance models, requirements, and data constraints on each other. Table 6 shows the dependencies of theoretical models, general definitions, data definitions, instance models, and likely changes on the assumptions.

[You will have to modify these tables for your problem. —TPLT]

[The traceability matrix is not generally symmetric. If GD1 uses A1, that means that GD1’s derivation or presentation requires invocation of A1. A1 does not use GD1. A1 is “used by” GD1. —TPLT]

[The traceability matrix is challenging to maintain manually. Please do your best. In the future tools (like Drasil) will make this much easier. —TPLT]

The purpose of the traceability graphs is also to provide easy references on what has to be additionally modified if a certain component is changed. The arrows in the graphs represent dependencies. The component at the tail of an arrow is depended on by the component at the head of that arrow. Therefore, if a component is changed, the components that it points to should also be changed. Figure ?? shows the dependencies of theoretical models, general definitions, data definitions, instance models, likely changes, and assumptions on each other. Figure ?? shows the dependencies of instance models, requirements, and data constraints on each other.

	TM??	TM??	TM??	GD1	GD??	DD1	DD??	DD??	DD??	IM1	IM??	IM??
TM??												
TM??			X									
TM??												
GD1												
GD??	X											
DD1				X								
DD??				X								
DD??												
DD??								X				
IM1					X	X	X				X	
IM??					X		X		X	X		
IM??		X										
IM??		X	X				X	X	X		X	

Table 4: Traceability Matrix Showing the Connections Between Items of Different Sections

9 Values of Auxiliary Constants

There are no auxiliary constants for this project.

	IM1	IM??	IM??	IM??	4.2.6	R??	R??
IM1		X				X	X
IM??	X			X		X	X
IM??						X	X
IM??		X				X	X
R??							
R??						X	
R??					X		
R2	X	X				X	X
R??	X						
R??		X					
R??			X				
R??				X			
R4			X	X			
R??		X					
R??		X					

Table 5: Traceability Matrix Showing the Connections Between Requirements and Instance Models

References

Hunter Ceranic. Software requirements specification for oar - optical alphabet recognition. 2024.

[The following is not part of the template, just some things to consider when filing in the template. —TPLT]

[Grammar, flow and L^AT_EX advice:

- For Mac users *.DS_Store should be in .gitignore
- L^AT_EX and formatting rules
 - Variables are italic, everything else not, includes subscripts ([link to document](#))
 - * [Conventions](#)
 - * Watch out for implied multiplication
 - Use BibTeX
 - Use cross-referencing
- Grammar and writing rules
 - Acronyms expanded on first usage (not just in table of acronyms)
 - “In order to” should be “to”

—TPLT]

[Advice on using the template:

- Difference between physical and software constraints
- Properties of a correct solution means *additional* properties, not a restating of the requirements (may be “not applicable” for your problem). If you have a table of output constraints, then these are properties of a correct solution.
- Assumptions have to be invoked somewhere
- “Referenced by” implies that there is an explicit reference
- Think of traceability matrix, list of assumption invocations and list of reference by fields as automatically generatable
- If you say the format of the output (plot, table etc), then your requirement could be more abstract

—TPLT]

	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??	A??
TM??	X																		
TM??																			
TM??																			
GD1		X																	
GD??			X	X	X	X													
DD1							X	X	X										
DD??			X	X						X									
DD??																			
DD??																			
IM1											X	X		X	X	X			X
IM??												X	X			X	X	X	
IM??														X					X
IM??													X					X	
LC??				X															
LC??								X											
LC??									X										
LC??											X								
LC??												X							
LC??															X				

Table 6: Traceability Matrix Showing the Connections Between Assumptions and Other Items