HORIZON-CL4-2021-DIGITAL-EMERGING-01
AI, Data and Robotics for the Green Deal (IA)

# AI-powered Robotic Material Recovery in a Box

## D6.4: Waste Data for material recognition and Recycling Data Game

| | |
|---|---|
| *Contractual Date of Delivery:* | 29/02/2024 |
| *Actual Date of Delivery:* | 15/04/2024 |
| *Security Class:* | Public |
| *Editor:* | *Antonios Liapis (UM)* |
| Contributors: | Dinos Larentzakis (ROBENSO), Freiderikos Raptopoulos (ROBENSO), Poojan Timilsina (IRIS), Laura Rodriguez (IRIS), Marina Pellegrino (IRIS), Michail Maniadakis (FORTH) |
| Quality Assurance: | Michail Maniadakis (FORTH) |
| Deliverable Status: | Final |

## The *RECLAIM* Consortium

| Part. No. | Short Name of Participant | Participant Organization name | Country |
|---|---|---|---|
| 1 | FORTH | Foundation for Research and Technology Hellas | EL |
| 2 | UoM | University of Malta | MT |
| 3 | KUL | Katholieke Universiteit Leuven | BE |
| 4 | HERRCO | Hellenic Recovery Recycling Corporation | EL |
| 5 | IRIS | Iris Technology Solutions, Sociedad Limitada | SP |
| 6 | RBNS | ROBENSO PC | EL |
| 7 | AIMPLAS | AIMPLAS - Technological Institute of Plastics | SP |
| 8 | AXIA | Axia Innovation UG | DE |
| 9 | ISWA | International Solid Waste Association | NL |
| 10 | ION | Periferiakos Foreas Diaxirisis Stereon Apovliton Ionion Nison Anonimi Eteria Ton Ota | EL |

## Document Revisions

| Version | Date | Editor | Overview |
|---|---|---|---|
| 0.1 | 16/02/2024 | Antonios Liapis (UM) | Structure and Early Draft |
| 0.2 | 28/02/2024 | Dinos Larentzakis (ROBENSO) | Details on Methodology and Dataset |
| 0.3 | 07/03/2024 | Antonios Liapis (UM) | Introduction and Future Work revisions |
| 0.4 | 15/03/2024 | Laura Rodriguez (IRIS) | Details about HSI data |
| 0.5 | 23/03/2024 | Antonios Liapis (UM) | Final edits and typesetting |

# Table of Contents

## List of Abbreviations

| Abbreviation | Definition |
|---|---|
| AI | Artificial Intelligence |
| DoA | Description of the Action |
| RDG | Recycling Data Game |
| prMRF | portable, robotic Material Recovery Facilities |
| RoReWo | Robotic Recycling Worker |
| HSI | Hyperspectral Imaging |

## Executive Summary

RECLAIM is a Horizon Europe funded project with an objective to develop a portable, robotic Material Recovery Facilities (MRFs) (prMRF) tailored to small-scale material recovery. RECLAIM adopts a modular multi-robot/multi-gripper approach for material recovery, based on low-cost Robotic Recycling Workers (RoReWos). An AI module combines imaging in the visual and infrared domain to identify, localize and categorize recyclables. The output of this module is used by a multi-RoReWo team that implements efficient and accurate material sorting.

Further, RECLAIM englobes a citizen science approach to increase social sensitivity to the Green Deal. This is accomplished via a novel Recycling Data-Game that enables and encourages citizens to participate in project RTD activities by providing annotations to be used in deep learning for the re-training of the AI module. Three different scenarios will attest its effectiveness and applicability in a broad range of locations that face material recovery challenges.

This deliverable provides an update on the dataset collection procedures in order to collect adequate waste data for two purposes: for training and testing AI algorithms and for soliciting ground truth data from citizen scientists through the recycling data game (RDG). Compared to D6.1 (submitted on M9), the updated dataset collection procedures are much closer to the realistic conditions of use of the prMRF, as all cameras, hardware, and technologies are already integrated in the prMRF and deployed on-site. This updated version is aligned with earlier planned future work reported on D6.1, and focuses on integrating with on-site developments with the prMRF and with the user data produced by the RDG (see D6.5 for the latter). The final pipeline for collecting data is expected to hold, with minor adjustments, throughout the lifetime of RECLAIM and the deployment of the prMRF on-location.

# 1. Introduction

This report covers the updated waste data collection procedures, which are expected to remain in effect throughout the lifetime of RECLAIM (with minor adjustments). Under the DoA, the goal of T6.1 is to *"allow the direct visual and hyperspectral examination of waste streams under the same conditions as the final operation of the prMRF after considering constant and efficient lighting on the waste."* The purpose of waste data collection within RECLAIM is two-fold. First, the AI algorithms for waste detection and categorization (under WP3) must be trained on both controllable and real waste data settings in order to assist the task of the robotic workers (RoReWo) in separating and recovering materials. Second, the waste data collected via processes described in this deliverable are used to produce recyclable data games (RDG) for the broader public to interact with (WP6). The gamified environment, described in D6.2, leverages citizen scientists to provide more – if less reliable – annotations than experts would be able to. Therefore, collecting a diverse set of images for the users to view while they play the game is vital; moreover, having control over such image datasets (e.g. knowing that a range of images contains only one specific type of object) allows for more gradual onboarding to players, leaving complex images with many different materials until later stages after players become familiar with the tasks (and materials) at hand.

Another goal of the methodology proposed in this report is, again according to the DoA, to *"be carried out periodically every 2 months for at least a whole year to meet the expected seasonal characteristics in the appearance of waste (e.g. raindrops in winter, dust in summer"*. Iterative waste data collection procedures are necessary in order to train more robust AI algorithms that do not learn to detect waste only on specific lighting or weather conditions. Importantly, iterative waste collection procedures will be paramount when the prMRF is deployed on-site as to refine current (largely controlled) experiments with the local context of the Ionian islands and the local communities, both in terms of light/weather conditions in the deployed prMRF but also regarding the composition of waste. Waste composition and volume can be very impactful to the quality of the AI algorithms, especially if those are trained only in controlled settings. Thus, iterative data collection once the prMRF is deployed will capture the local context and conditions (e.g. in the summer, more dust and larger volumes of waste due to tourist influx, versus darker lighting conditions in the fall and less volume from the local population). Based on first experiments facilitated from the controlled dataset collected from D6.1 (submitted M9), feasibility tests and development of the RDG (with D6.5 submitted concurrently to this report) allowed us to refine use-cases regarding how data can be useful for engaging the public in annotation tasks. With ongoing prMRF development (see D5.2 submitted concurrently), the updated methodology for collecting data is more ecologically valid and expected to hold, as a method, for collecting on-location data in a realistic and sustainable fashion for both AI and RDG purposes.

## 1.1  Intended readership

The present report is a public (PU) document. Its readership is considered to be the European Commission, the RECLAIM Project Officer, the partners involved in the RECLAIM Consortium, beneficiaries of other European funded projects, and the general public.

## 1.2   Relationship with other RECLAIM deliverables

The methodology and collected data described in this report have been developed in light of the data management plan and ethics/privacy manual reported as part of D1.1. The design of the waste data collection methodology, and format of collected data, was informed by the needs of the AI algorithms for material identification, localization and categorization under WP3 (T3.1, T3.2) and the needs of the recycling data games that will make such data public and accessible to users (T6.2 and T6.3). Table 1 shows the main deliverables consulted (in case of past work), and impacted by (in case of future work) by this report.

| Del. No | Deliverable Name | WP | Month |
|---------|------------------|----|-------|
| 1.1 | Data management plan and ethics/privacy manual | WP 1 | M6/M36 |
| 2.1 | prMRF and RDG requirements and systems specification | WP 2 | M6 |
| 3.1 | Material recognition based on RGB and Hyperspectral imaging | WP 3 | M18 |
| 3.2 | prMRF operation monitoring and repeating advancement | WP 3 | M30 |
| 6.1 | Waste Data for material recognition and Recycling Data Game | WP 6 | M9/M18 |
| 6.2 | Algorithms and pipelines for Recycling Data Games | WP 6 | M9/M18/M30 |
| 6.3 | Assessment of the Recycling Data Game | WP 6 | M18/M36 |
| 1.3 | Intermediate / Final Project Report | WP 1 | M18/M36 |

*Table 1: Other RECLAIM deliverables related.*

## 2. Updated waste data collection methods: Integration of waste data collection within the prMRF

As discussed in Section 1, the updated methodology for collecting waste is integrated with the work-in-progress of the prMRF (portable robotic Material Recovery Facility). Thus, the data collection is more ecologically valid and expected to hold for collecting on-location data in a realistic and sustainable fashion, even when the prMRF is moved during the needs of the RECLAIM project, and afterwards.

The current build of the prMRF houses three camera systems, complemented by illuminating lights, two sets of linear robots, and a magnetic separator and vibratory table. Below, we describe the sequence of data capturing within the prMRF (see Section 2.1) and provide more details about the illumination and data capturing devices (see Section 2.2). We conclude this section with the way we collected the revised dataset (see Section 4) using the prMRF itself.

### 2.1 prMRF waste trajectory and data collection pipeline

The prMRF setup includes three data collection steps: one HSI (HyperSpectral Imaging) with associated illumination setups, and two RGB camera boxes which include illumination devices. Below we describe their location in the prMRF and, implicitly, the timing of the data capture in the process of the waste processing chain handled by the prMRF.

The prMRF waste processing pipeline starts with the vibratory table for spreading out the recycling material on the conveyor belt, which operates at a speed ranging from 200 to 500 mm/sec (in RECLAIM, a typical conveyor belt speed is expected to be around 300mm/sec).

Right after that, the **Visum HSI™ hyperspectral imaging system** adapted for the RECLAIM application is installed, coupled with the corresponding illumination units. The HSI camera is installed at a height of 60cm from the conveyor belt, scanning vertically to the flow of the moving belt. The lighting module consists of 2 series of halogen lamps, each one assembled in a steel frame which hangs from the prMRF ceiling. The two lighting modules are installed before and after the Visum system with a certain angle focused on the location that the HSI camera is line-scanning (see Figure 1).
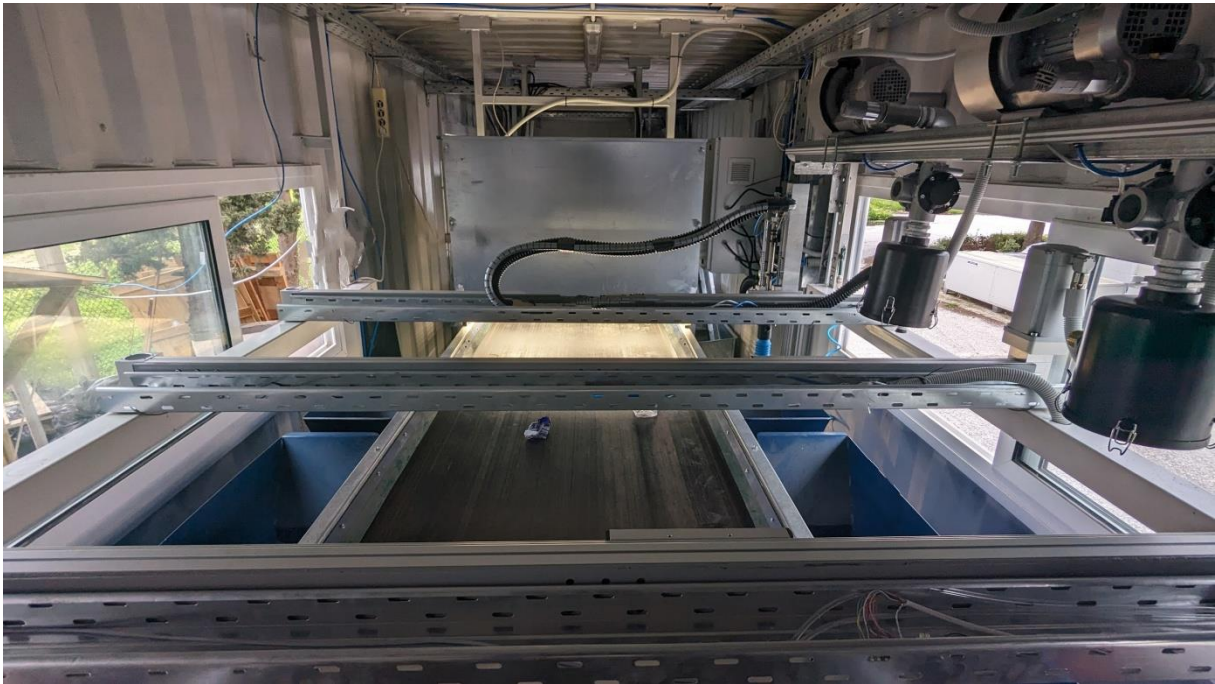
*Fig. 1: The adapted Visum HSI™ hyperspectral imaging system with its lights mounted front and back (side view)*

Following this, the first set of RGB camera boxes is installed (see Figure 2). Each box consists of (i) a dark room to to prevent interference and minimize the effect of external lighting sources, (ii) a 2-megapixel RGB camera, housed within a metal frame suspended from the prMRF's ceiling on a steel base, (iii) an internal lighting arrangement that consists of twelve white illumination LED bars. Eight of these bars are evenly distributed and horizontally aligned at the same height as the camera (90 cm from the conveyor belt). The remaining bars are grouped in pairs, positioned at the left and right side of the camera (10cm lower than the camera) in order to enhance uniform lighting along the entire length of the conveyor belt.

Subsequently, the first team of Robotic Recycling Workers (RoReWos) is installed. It consists of three 1.5 DOF linear robots (see Figure 2). The main responsibility of this first section, HSI and RGB cameras and the RoReWo team, is to identify and retrieve materials belonging to a predefined subset of the following six classes of recyclables: PET, HDPE, Aluminum, Tetrapak, PP/PS, and PE Film. Additionally, a magnetic separator efficiently collects ferrous materials, diverting them along a separate disposal path.

The waste processing chain proceeds with the second set of RGB camera boxes housed in an identical metal enclosure with an identical lighting setup and camera as the one described above, before concluding with the second RoReWo team, featuring one 3.0 DOF and one with 2.5 DOF linear robots (see Figure 3). The second RoReWo team is targeted in identifying and locating materials not captured in the initial sorting phase, encompassing the remaining volume of the initial six classes and additional the remaining classes of materials PE Films and Other plastics.

*Fig. 2: The first section with the HSI camera (behind the metal box not included) and the first RGB camera box and the 3 1,5 DOF linear in series.*



*Fig. 3: The second section after the magnet of the prMRF with RGB camera box with lights and the two linear 3 DOF (first from the camera box) and the 2.5 DOF (the second one).*

## 2.2 Lighting and camera setup in the prMRF

The lights inside the darkroom of the RGB camera that is made of a steel frame consist of 12 LED bars of white illumination (see Figure 4). These lights provide a moderate level of brightness to prevent reflections on the recycling materials. They are arranged horizontally in a row at the same height as the RGB camera, which is 90 cm, and mounted within a steel construction. On both the left and right sides of the camera, there are 4 bars spaced 6 cm apart. Additionally, at each corner, there are 2 bars positioned 10 cm lower than the others to ensure better illumination while avoiding reflections from the side panels of the steel box.

The illumination unit for the HSI camera consists of 2 series of 6 x 50Watt halogen lamps of 12Volts each that are placed in front and back of the camera on 25cm height of the conveyor belt. Calibration unit is also embedded in the system used to correct the final results.



*Fig. 4: The setup for the RGB camera with the light bars.*

The RGB cameras (one in each camerabox) are Baumer 2-megapixel RGB cameras with a full frame 1920X1200 pixels with 59 max fps (frames per second). The speed of data collection that we use inside the prMRF is 10 frames per second (fps).

The hyperspectral imaging system consists of the hyperspectral sensor, the illumination unit, the calibration unit, and the embedded computer that analyzes the data. The current system is capturing 100 lines per second. The system directly exports the prediction of every pixel based on the seven categories of the specified material (PET, HDPE, PP/PS, Aluminum, Tetrapak, PE Film). The final result is sent to the other systems via the Mosquito broker and the MQTT protocol [1], serving as a communication layer with low latency.

## 2.3 Current dataset collection methodology

The initial data collection technique summarized in D6.1(M9) involved mixed flow pathways of recyclables captured by cameras placed within the recycling center of Heraklion (operated by HERRCO), and industrial sites where ROBENSO had installed computer vision modules. This was mainly because the RECLAIM prMRF (portable robotic Material Recovery Facility) was not yet developed.

Due to the high density of materials in these images, the annotations of recyclable objects were mostly done manually, as it was challenging even for state-of-the-art tools like Segment Anything (see D6.5) to easily identify the borders separating recyclable materials, especially in the case of (semi-) transparent objects. This has been a very time-consuming process, as annotating a single image took approximately 2.5 minutes. The latter makes the acquisition of large and well-annotated image datasets particularly challenging. Another issue that appeared with the manual annotation of the industrial images is the uneven distribution of materials in the flow, resulting in significant disparities in data quantity for each material. This means that, indicatively, there was a need to annotate thousands of PET but only a few hundreds of HDPE objects. This imbalance could significantly affect the neural network's recognition efficiency, especially for less frequent materials.

To overcome the above issues, we have adopted a new methodology that is more straightforward and does not require extensive manual labor. In particular, we focused on utilizing pre-sorted materials provided by HERRCO. These materials are received from the collection bins, where HERCCO workers place the sorted recyclables. As a result, we have gained access to a significant number of objects per material type, allowing us to have separate material streams for each category we aim to identify. Following this approach, we have been able to semi-automate annotation, as all identified objects of a given stream are classified in the same category. What remained to be solved was the accurate segmentation of objects. By making sparse streams, we have been able to make quicker annotation, approximately half the time compared to the previous method. This is because segmentation could be done largely automatically by our own trained neural networks or the Segment Anything tool, with subsequent small scale manual adjustments/corrections. Interestingly, the

annotated data and the detailed object masks have been utilized to create synthetic images, potentially being able to multiply our data for a given recyclable material by tenfold.

Additionally, we have also examined the use of mixed streams but still with controlled flows and material density, making slight adaptations on the prMRF conveyor belts. The data obtained this way provide the basis for generating  particularly realistic (yet synthetic) images of mixed materials. This is achieved by using the masks discussed above to generate multiple realistic synthetic single object images, particularly for sparser materials, which are used to significantly increase our dataset.

In both approaches discussed above, random physical placement of recyclables on the conveyor belt is preferred using the already available looping conveyors operating in the prMRF. However, manual intervention was occasionally necessary for either creating challenging conditions where materials overlap, or for better arrangement of materials to obtain clear masks for synthetic data generation. Due to the use of global shutter RGB cameras which capture clear images without motion distortions, the speed of the conveyor belt has no effect on the collected data. Still, the speed of the conveyor belt can be used to adjust the spread of the materials. To collect the data summarized above we used speeds that occasionally varied from half to the full targeted belt speed within the container (i.e. from 150-300 mm/sec). The use of slower conveyor belt speeds allowed the ROBENSO staff to adjust the stream (arranging of materials) as needed.

By employing this methodology, we streamlined the data collection process while ensuring a diverse and comprehensive dataset for training neural networks that actually implement the recyclable identification, localization and categorization module. Currently, our dataset comprises 1200 images for each individual class material and 5000 images of mixed streams (see Section 4), providing a robust foundation for training and testing algorithms in diverse operational scenarios. This data will be made available for public use (see Section 5).

## 3. Data Format

Given the extensive collaborative decision-making undertaken for the purposes of D6.1, we largely retained the same data format for waste data. While "live" data during prMRF deployment will be added to the database (described at the end of this section), for the sake of posterity and backup, we continue to keep data in Google Drive in folders that allow researchers to track the provenance of the data. At the root folder, folders will denote the starting date of data collection for this batch (for instance, if data collection starts on 1 May 2023, the folder will be named 2023.05.01). Within each of these folders, there will be subfolders denoting the type of data (RGB versus HYPER) and within each of those there will be folders with the stream used to produce the data (e.g. "mixed", "PET", "ALU"). A consistent naming convention will be retained for these folders that matches the remaining data format naming process within RECLAIM (e.g. as outputs of the AI algorithms). These folders will contain all data pertaining to this data collection period and process, for example including both the image and the json files that describe it, using the same file naming convention (and when needed, consistent suffixes).

RGB images are stored in the JPEG format, at resolution 1920x1200 pixels. Such images are on the file size scale of 500 kilobytes. Each dataset containing a certain group of images is accompanied by a JSON file that describes the metadata of the dataset. The latter includes the information content of each image as predicted by the current AI algorithms implemented in WP3. The JSON file contains information about our annotated pictures, divided into four sections. The first section provides details such as the date and file name. In the second section, we find the resolutions, name of the photo, and unique IDs assigned to each photo, whether or not it has been annotated. Moving on to the third section, it contains all the annotations created, including the annotation ID, photo ID, the number of points defining the height and width of each point that we made for every annotated object, and the assigned category ID. Lastly, the fourth section lists the material categories along with their corresponding IDs that we have established.

Concurrently, we have developed a database for the purposes of RDG, which can at the time of writing store RGB and metadata (as JSON files described above). D6.5 contains details about the database and querying mechanisms. The database is live, and will be populated by current collected data produced via the pipeline mentioned in Section 2.3 (and currently on Google Drive). In the coming months (see Section 5) new data produced via the process described in Section 2 within the prMRF will be automatically stored in the RDG database. Additional development for the RDG database, for example to include HSI data, will be investigated based on both technical feasibility and user's feedback under T6.2: Algorithms and pipelines for recycling data-games and T6.3: Assessment of the Recycling Data Game.

## 4. Current dataset

| Data type | Waste type | Size |
|---|---|---|
| RGB images | PET | 1200 images |
| RGB images | HDPE | 1200  images |
| RGB images | PE Film | 1200 images |
| RGB images | Aluminum | 1200 images |
| RGB images | Tetrapak | 1200 images |
| RGB images | PP/PS | 1200 images |
| RGB images | Mixed | 5000 images |
| Hyperspectral data | PP | 1780 hyperspectral images |
| Hyperspectral data | PS | 900 hyperspectral images |
| Hyperspectral data | PET Bottles | 3000 hyperspectral images |
| Hyperspectral data | Tetrapak/Paper | 1000 hyperspectral images |
| Hyperspectral data | HDPE | 1500 hyperspectral images |
| Hyperspectral data | LDPE | 1000 hyperspectral images |
| Hyperspectral data | Aluminum cans | 1000 hyperspectral images |
| **Total** | **12200 RGB images, 10180 hyperspectral images** | |

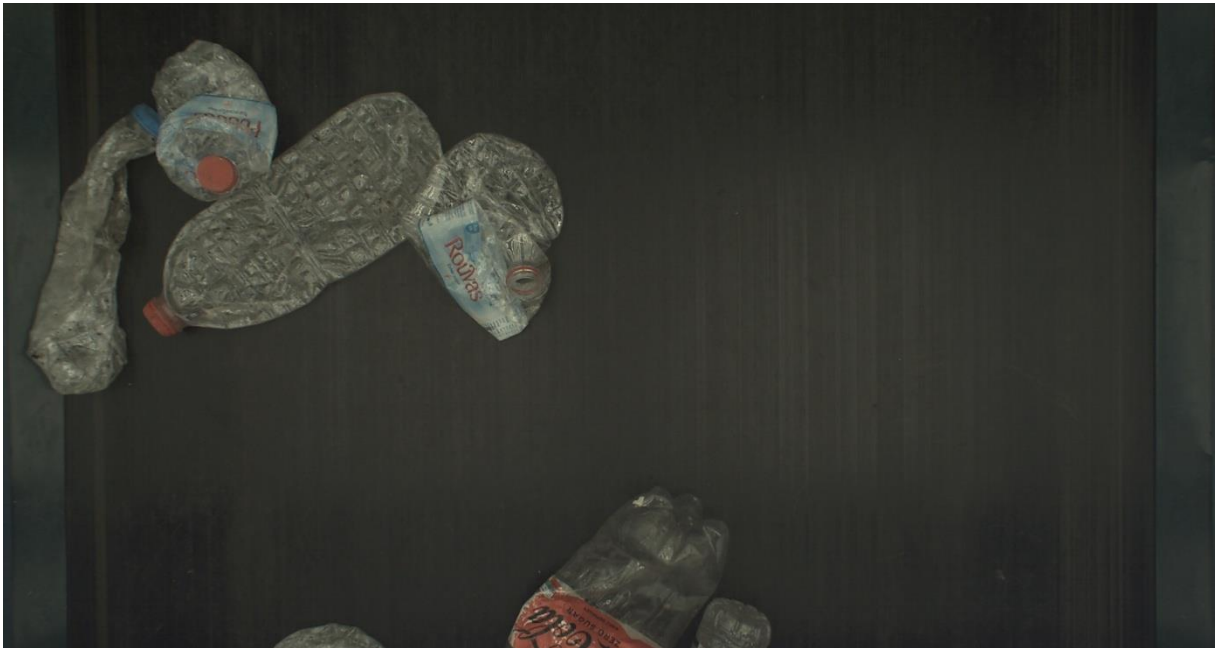*Table 2: Contents of the current dataset split per data type and material.*

HSI images acquired so far for the current dataset are summarized in Table above. Hyperspectral data are stored in ENVI files. The ENVI file format is standardized. The ENVI file contains the collection of HSI data.

To be used in this dataset, the updated size of every image is 1 x 224 x 640. All the images selected to be included in the current dataset are stored in 38 ENVI image collection files. Along with the enlargement of the dataset with new HSI images in the future, the format of the HSI data in the dataset will be improved as well.
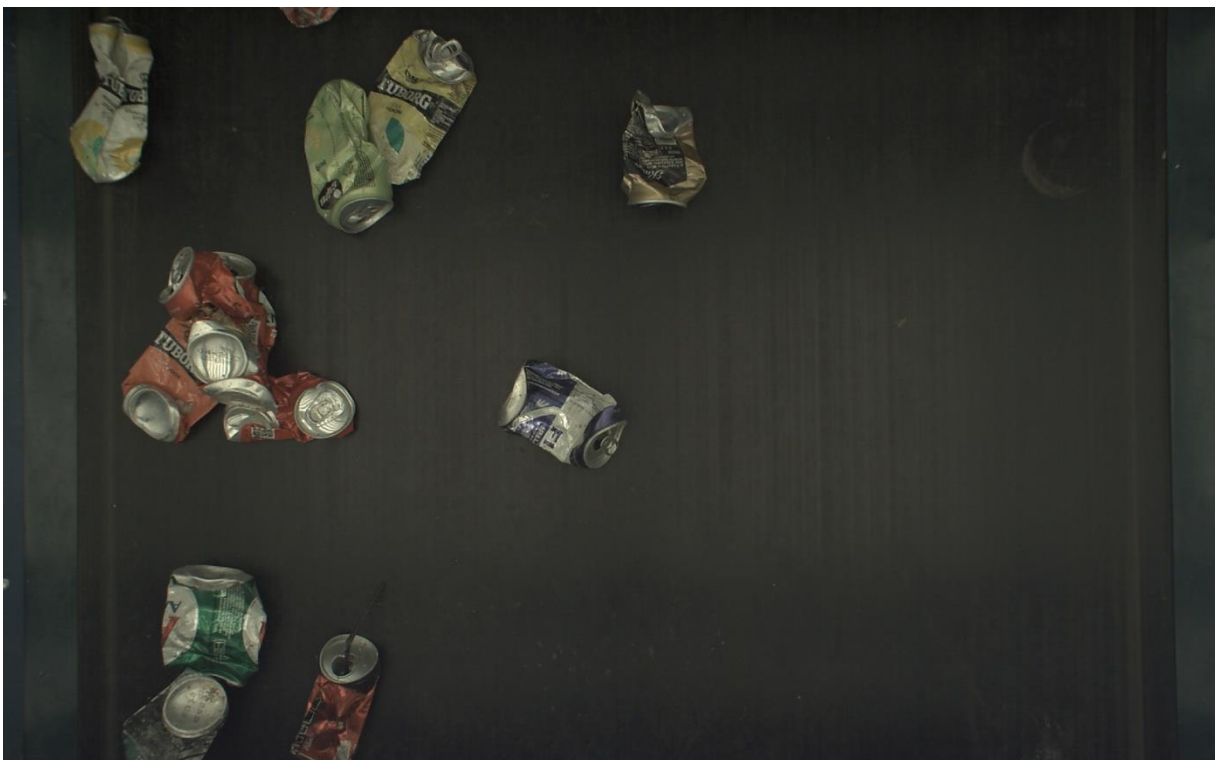
Table 2 summarizes the dataset collected so far based on the methodologies described in Section 2.3. At the current point in time, we have collected 12,200 RGB images (with the possibility to collect many more via the mixed recyclable stream described in Section 2.1.1) as well as 10,180 hyperspectral images. The hyperspectral images collection will continue to
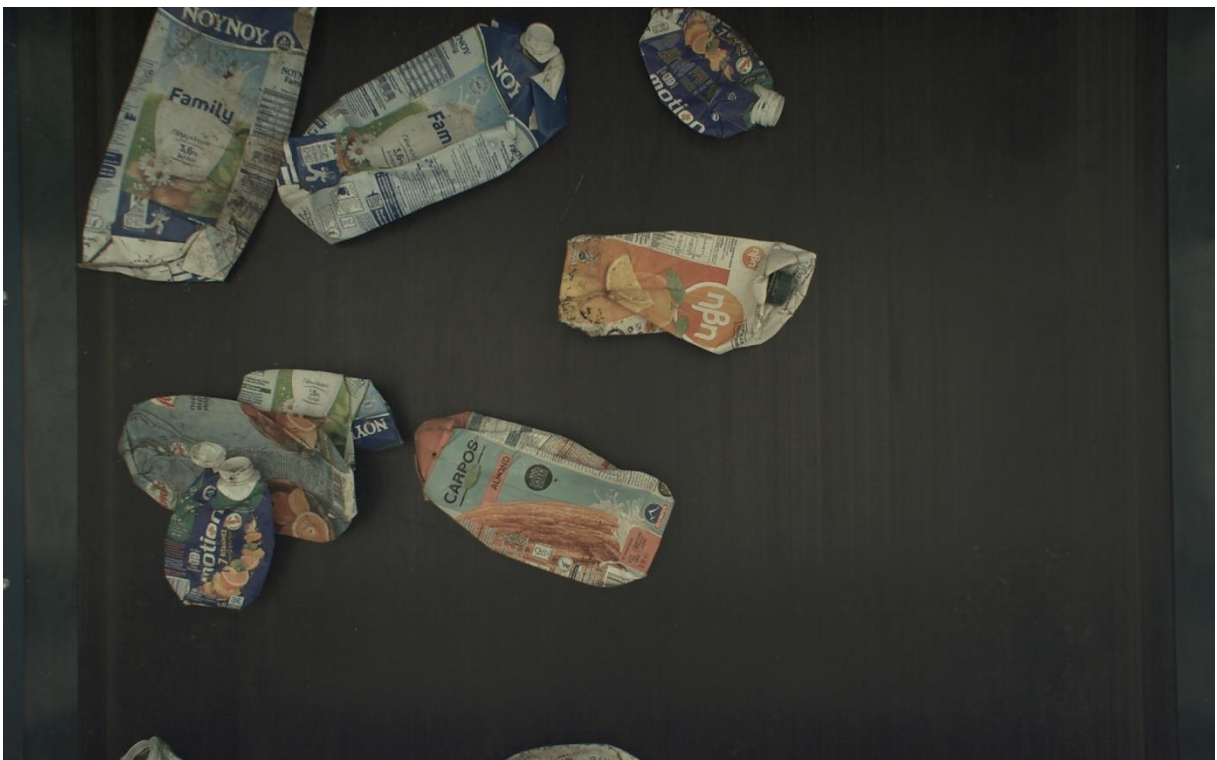
enlarge the number of images available. This offers a very rich dataset which can allow us to kickstart training of AI algorithms as well as to present to players of RDGs.

Below we offer a snapshot of the data collected in the above dataset. Figures 5-9 show the RGB images captured from the isolated streams, with different materials. Figure 10 shows an RGB image from a controlled mixed waste stream (with specific materials) overlaid with AI annotations of different materials as polygons. The data for these AI annotations is stored in a JSON file and the visualization can be adjusted as needed (e.g. for different resolutions in the Recycling Data Game). D6.5 already includes mini-games where these AI annotations are shown to the users (and curated by them). We note that the images in this dataset are different from D6.1 in important ways, as the illumination (which has been fine tuned in the intervening months in the actual prMRF) minimizes reflection that may challenge AI algorithms for detecting waste (see WP3).



*Fig. 5: RGB image: Mixed waste stream.*

*Fig. 6: RGB image: Isolated PET waste stream.*



*Fig. 7: RGB image: Isolated Aluminum waste stream.*

*Fig. 8: RGB image: Isolated PP_PS waste stream.*



*Fig. 9:  RGB image: Isolated Tetrapak waste stream.*

*Fig. 10: Annotated RGB image: mixed waste stream with PET, Aluminum, Tetrapak and Other objects annotated.*

# 5. Future Work

This deliverable concludes T6.1, with a methodology for collecting waste data that is both ecologically valid (data collected is from within the prMRF) and sustainable (as data can be stored directly in the RDG database prepared under D6.5). However, several issues remain for future work including deployment of the prMRF, collection of more data under "real" conditions, and testing the RDG database using live data.

The method followed for producing the dataset in Section 4 is inspired by earlier versions of this methodology (D6.1) where controllability is key. As the data collection process has been largely finalized (see Section 2), it is expected that the collection of new waste data from live deployment will largely lead to similar results as those in our controlled experiments reported in Section 4. However, technical adjustments may be required when the prMRF is on-site, and minor adjustments to the methodology may be required (including, for example, the speed of the conveyor belt for both capturing crisp images and for actual material sorting).

The dataset currently is on shared storage spaces (Google Drive) for the purposes of tracking its provenance and backup. However, the RDG database is already in place (D6.5) and already contains a subset of the dataset of D6.1. Current experiments on the RDG (see D6.3 and D6.5) seem to handle the data format well, but additional experiments will need to be conducted. The first experiment is regarding the "live" feedback loop from the RGB cameras and AI algorithms to the database, in order to assess the speed of data input and the resulting volume of data that needs to be shown to players. The second experiment is regarding the applicability of the "live" data for the RDGs. Minor adjustments (in terms of resolution or brightness levels) is expected on the game side – without affecting the original image. Finally, user data is intended to be summarized as a new ground truth for the AI algorithms to leverage; while development is already underway for this in D6.5, integration between WP3 and WP6 is still expected to be needed to harmonize the feedback loop.

In the upcoming months, XXX will primarily focus on the data collected from citizens through RDG playing. Efforts will be made to make the dataset available to the public, in the principle of FAIR [2] and open science. Once the dataset above is cleaned, documented, and verified in terms of non-traceable data, it will be placed in an open science repository (to be determined and published in the RECLAIM website) and made available for AI training or other uses.

# 6. Conclusion

This report presented the current methodology implemented by ROBENSO in a working version of the prMRF, in collaboration with IRIS, FORTH and University of Malta, for collecting datasets of waste under different conditions. While the methodology is now ecologically valid, as it represents the real-world conditions (including lighting), the data collected remains controllable following the methodology initially presented in D6.1 and refined as discussed in Section 2.3. Through this methodology, a dataset of 12200 RGB images and 10180 HSI images has already been collected, with more data acquisition planned for both RGB and HSI images. The dataset is made available as part of D6.4. The method will remain in effect during on-location experiments in data collection, and refined as needed to maximize the possible data shared with the RDG (and therefore, the public). Moreover, work from now on will focus on collecting diverse waste material data via this sustainable method on different lighting and weather conditions.

# References

[1] R. A. Light, "Mosquitto: server and client implementation of the MQTT protocol," The Journal of Open Source Software, vol. 2, no. 13, May 2017, DOI: 10.21105/joss.00265

[2] Wilkinson, M., Dumontier, M., Aalbersberg, I. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **3**, 160018 (2016). https://doi.org/10.1038/sdata.2016.18