

# STAIRWAI

*Stairway to AI: Ease the Engagement of Low-Tech users to the AI-on-Demand platform through AI, H2020*

## Benchmarking software-1st version

Deliverable information	
Deliverable number	D6.1
WP number and title	WP6 - Vertical Matchmaking and hardware marketplace
Lead beneficiary	BCA
Dissemination level	Public
Due date	31 August 2022
Actual date of delivery	19 October 2022
Author(s)	Miguel de Prado (BCA), Jean Marc Bonnefous (BCA)



## Document Control Sheet

Version	Date	Summary of changes	Author(s)
0.1	15/07/2022	Initial draft of document	Miguel de Prado (BCA)
0.2	22/07/2022	First version of deliverable completed	Miguel de Prado (BCA)
0.3	29/08/2022	Internal review completed	Michela Milano (UNIBO), Gabriel Gonzalez (Insight)
0.4	23/09/2022	Answers to reviewer's comments added	Miguel de Prado (BCA)
1.0	19/10/2022	Final version of the document	



## Table of contents

1. Executive Summary	5
2. Introduction	6
2.1. Purpose and Scope of the document	6
3. LPDNN	7
3.1. LPDNN architecture	7
3.2. Inference engines	8
3.3. Supported HW platforms	9
3.4. Benchmark framework	10
4. Deployment of AI Apps and Benchmark	11
4.1. Deployment of AI Apps	11
4.1.1 LPDNN AI Application (AI App)	11
4.1.2 LPDNN Deployment Package	12
4.2. Benchmark of AI Apps	12
4.2.1. Choose AI App and target HW platform	12
4.2.2. Requirements	13
4.2.3. Download AI App and deployment package	13
4.2.4. Benchmark your AI-App on your target HW	14
4.3. Benchmark as a Service	15
5. Conclusion and future work	17
Bibliography	17



## Acronyms

Acronym	Explanation
AI	Artificial Intelligence
API	Application Programming Interface
BMP	Bonseyes Marketplace
DNN	Deep Neural Network
LPDNN	Low-power Deep Neural Network framework
ML	Machine Learning
WP	Work Package



# 1. Executive Summary

This document delivers the first version of the *Benchmark software framework* proposed within WP6 Task 6.1 of StairwAI project. The *benchmark software framework* represents one of the main blocks of the vertical matchmaking service as it produces benchmarks across several HW platforms (Task 6.2) that are then used to train the vertical matchmaking engine (Task 6.3).

Section 2 provides a general introduction of the WP and summarize the scope and contributions of this deliverable.

Section 3 introduces LPDNN: the inference framework that has been used and extended to create a benchmark framework for heterogeneous HW platforms. First, LPDNN's architecture is detailed, explaining the different components that make it an interoperable framework. Next, we introduce the different inference engines that LPDNN supports, which increase the portability and optimisation of neural networks across HW platforms. Then, we describe the support for HW platforms and those platforms that have been already integrated.

Section 4 starts by explaining the elements that are used to deploy AI Applications on HW platforms. Then, the benchmark flow is explained by showcasing the execution of an AI Application and the obtention of a benchmark. Finally, the creation of a Benchmark as a Service and its integration into AI4EU is discussed.

Lastly, Section 5 elaborates the conclusions and future work on Task 6.1.



## 2. Introduction

WP6 has the main objective of building a Vertical Matchmaking engine that matches AI algorithms and HW resources to optimise the deployment of services and increase their efficiency.

As such, one of the main objectives of WP6 is to develop *benchmark software framework/service* suitable for selected machine learning models and hardware, including CPU (Intel x86, Arm Cortex-A5x, Risc-V), GPGPU (NVIDIA, etc), and NPU (HUA, etc) accelerated platforms. The benchmark software framework will be used in Task 6.2 to produce the benchmarks and create a profiling dataset that can be used in Task 6.3 to train the Vertical Matchmaking engine's algorithms.

In this document, we propose LPDNN framework as the *benchmark software framework*. LPDNN framework was developed during the H2020 Bonseyes project (Prado, Miguel De, et al) and has been largely extended to provide a structured benchmarking layer to analyze the execution of AI Applications on a variety of heterogeneous platforms. This layer, on top of LPDNN, adheres to the following design principles:

- ❑ Industry-driven Research:
  - Input requirements from SMEs wanting to use AI solutions
  - Able to deploy AI solutions on edge devices (low-power, low-carbon footprint)
- ❑ Structured benchmarking workflow:
  - Availability of anchors within the deployment framework to evaluate the metrics truthfully
  - Optimised deployment (value added to research and industry)
  - Extensive documentation & Support (user friendly for SMEs)
  - Defined interfaces (standardization)
  - Easy to replicate (reproducibility)
  - Create trust & community

Taking the previous design principles, the main contributions of this deliverable are the following:

- Introduction of LPDNN as an inference framework.
- Integration into LPDNN of three inference engines to benchmark CPU, GPU and NPU platforms.
- Integration into LPDNN of three HW platforms and update of four available platforms to the latest SW release.
- Creation of the benchmarking layer within LPDNN, containing a variety of static and dynamic metrics.
- Evaluation of the inference/benchmarking framework as Service that can be used by external users and integrated into the AI-on-demand platform.

### 2.1. Purpose and Scope of the document

Deliverable D6.1 is a Demonstrator, i.e., it introduces the first version of the *Benchmark software framework*, which will be extended during the course of Task 6.1 and will conclude with D6.2.



### 3. LPDNN

LPDNN, which stands for Low-Power Deep Neural Network, is a deployment framework that provides the tools and capabilities to generate portable and efficient implementations of DNNs. The main goal of LPDNN is to provide a set of AI applications for deep learning tasks, e.g., object detection, image classification, speech recognition, which can be deployed and optimised across heterogeneous platforms, e.g., CPU, GPU, FPGA, NPU (ASIC).



Figure 1 LPDNN AI classes

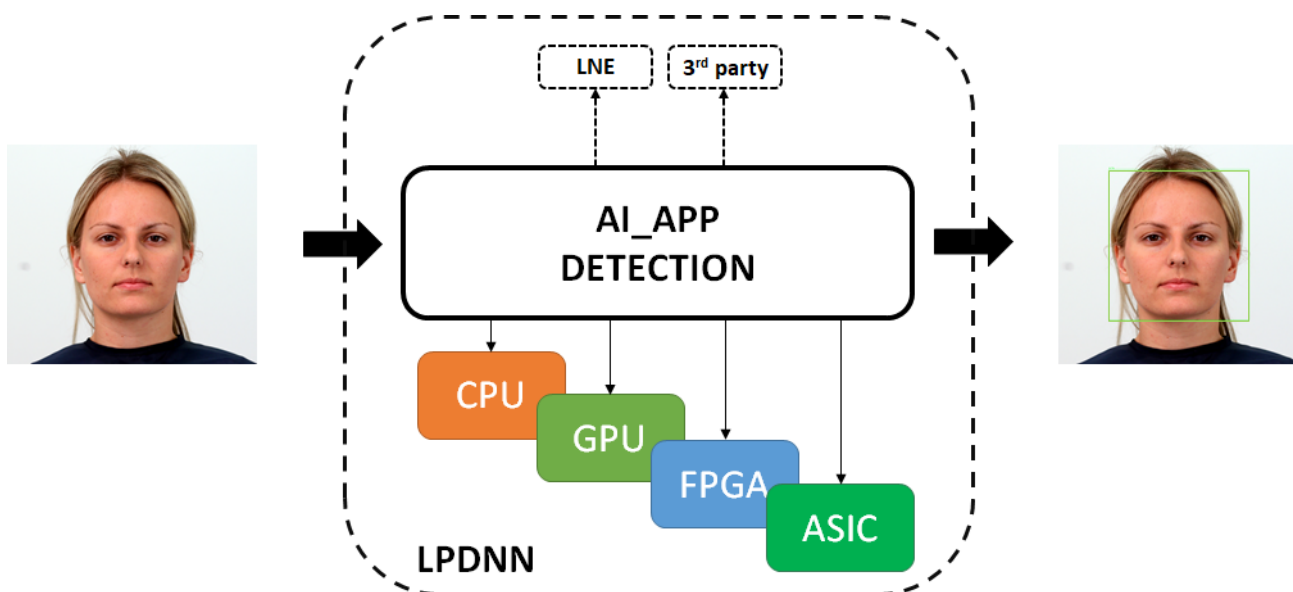


Figure 2 LPDNN overview

#### 3.1. LPDNN architecture

One of the main issues of deep learning systems is the hardship to replicate results across different systems. To solve this issue, LPDNN features a full development flow for deep learning solutions on embedded devices by providing platform support, sample models, optimisation tools, integration of external libraries and benchmarking. LPDNN's full development flow makes the AI solution reliable and easy to replicate across systems.



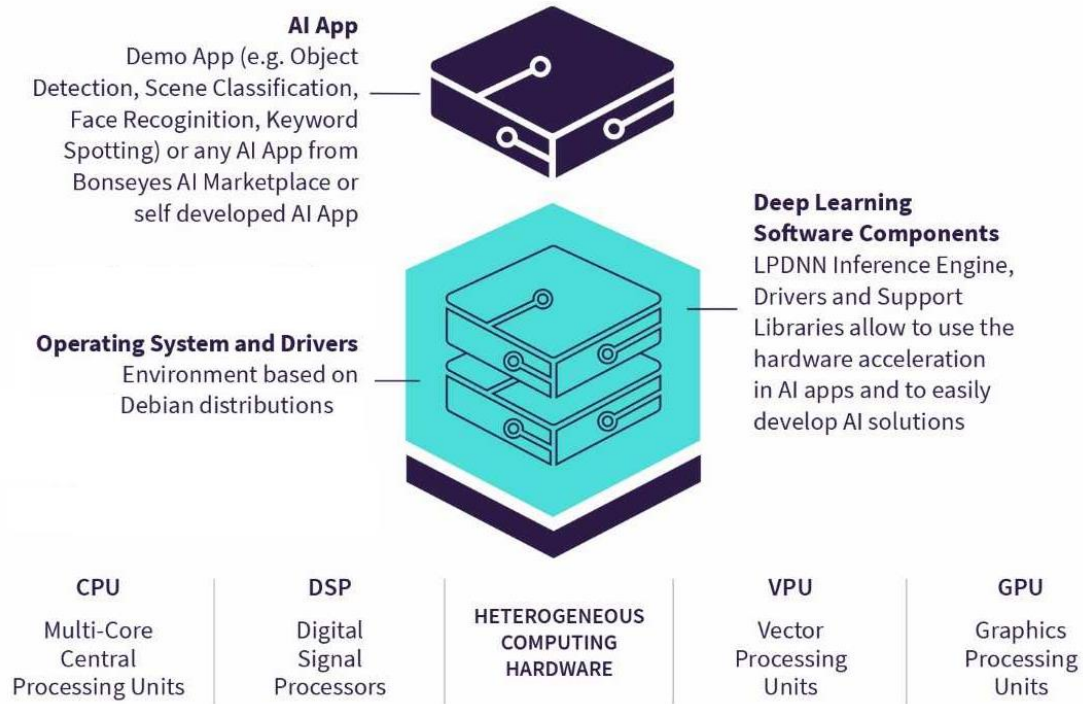


Figure 3 LPDNN full stack

AI applications are the result of LPDNN’s optimisation process and the higher level of abstraction for the deployment of DNNs on a target platform. They contain all the necessary elements or modules for the execution of a DNN. An AI application is the optimised representation of a DNN model to be efficiently executed on a target embedded device. An AI application contains all the necessary elements or modules for the execution of a DNN:

- **Pre-processing:** Step to prepare, normalize or convert the input data into the required input that is expected by the DNN.
- **DNN inference:** Forward-pass of the neural network. The execution is taken care of by an inference engine.
- **Post-processing:** Conversion of the neural network’s output into structured and human-readable information.

Next, we detail LPDNN’s architecture by further describing the concept of LPDNN ‘s inference engines and the support for heterogeneous platforms.

### 3.2. Inference engines

AI applications contain a hierarchical but flexible architecture that allows new modules to be integrated within the LPDNN framework through an extendable and straightforward API. For instance, LPDNN supports the integration of 3rd-party self-contained inference engines to perform DNN inference. Initially, LPDNN only supported:

- **LNE:** [LPDNN Native Engine](#) (LNE) allows the execution of DNNs across arm-based and x86 CPUs as well as on Nvidia-based GPUs.

During the development of the Stairwai project, the following inference engines have been integrated:



This project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 101017142



- **NCNN**: [NCNN](#) ports the execution of DNNs on GPU through the Vulkan driver.
- **TensorRT**: [TensorRT](#) accelerates the DNN inference on Nvidia-based GPUs and NPUs.
- **ONNXruntime**: [ONNXruntime](#) enables the direct execution of ONNX models on CPUs and GPUs.

The inclusion of external engines also benefits LPDNN as certain embedded platforms provide their own specific and optimised framework to deploy DNNs on their hardware.

More information about how to add a new inference engine in LPDNN can be found at:

- Developer guides: [https://bonseyes.gitlab.io/bonseyes-cli/pages/dev\\_guides/ai\\_app\\_index.html](https://bonseyes.gitlab.io/bonseyes-cli/pages/dev_guides/ai_app_index.html)

### 3.3. Supported HW platforms

One of the main factors for LPDNN's adoption is performance **portability across the wide span of hardware platforms**. LPDNN's flexible architecture allows the main core to remain small and dependency-free while additional 3rd party libraries or inference engines are only included when needed and for specific platforms, notably increasing the portability across systems. Besides, cross-compilation and specific tools are added to support a **wide range of heterogeneous computing platforms such as CPUs, GPUs, NPUs**. One of the objectives of LPDNN is to provide full support for reference platforms by providing:

- **Developer Platform Environments (DPEs)** to help the user employ a developer platform, including OS images, drivers, and cross-compilation toolchains for several heterogeneous platforms.
- **A dockerised and stable environment**, which increases the reliability by encouraging the replication of results across platforms and environments.
- **Optimisation tools and computing libraries** for a variety of computing embedded platforms that can be used by LPDNN's inference engines to accelerate the execution of neural networks.

Originally, the range of embedded platforms that were supported within LPDNN were the following:

- **Raspberry Pi 3b+**: Quad-core ARM Cortex-A53 (ARMv8) 64-bit SoC @ 1.4GHz
- **Raspberry Pi 4b**: Quad-core ARM Cortex-A72 (ARM v8) 64-bit SoC @ 1.5GHz
- **Nvidia Jetson Nano**: Quad-core ARM Cortex -A57 @ 1.43 GHz & 128-core Nvidia Maxwell GPU
- **Nvidia Jetson Xavier**: Octa-core ARM v8.2 @ 2.03 GHz & 512-core Nvidia Volta GPU with Tensor Cores

During the development of the Stairwai project, we have employed Bonseyes LPDNN's platform workflows to integrate new HW platforms and provide a more heterogeneous set of benchmarks on variety of HW processing cores. The following HW platforms have been fully integrated:

- **Intel NUC**: Intel quad-core (TM) i5-7260U CPU @ 3.4 GHz
- **iMX8m Nano**: Quad-core ARM Cortex-A53 (ARMv8) 64-bit SoC @ 1.4GHz
- **STM32 MP1**: Dual-core ARM Cortex-A7 cores up to @ 800 MHz

Besides, the following platforms have been upgraded:

- **Nvidia Jetson Nano**: Update to latest jetpack release JP6.4
- **Nvidia Jetson Xavier**: Update to latest jetpack release JP6.4
- **Raspberry Pi 3b+**: Update to latest Ubuntu20 packages



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 101017142

- **Raspberry Pi 4b:** Update to latest Ubuntu20 packages

Furthermore, BonsAPPs project is currently performing the integration of extra-low-power platforms containing Micro-controller Units, which can then be leveraged in Stairwai for benchmarking:

- **STM32 H7A3:** Arm® Cortex®-M7 core @ 480 MHz & 1.4 Mbytes of SRAM
- **Greenwaves GAP8:** RISC-V core & Octa-core RISC-V (cluster) & 512kB of L2 memory

More information about the HW platforms can be found at:

- User guide: [https://bonseyes.gitlab.io/bonseyes-cli/pages/user\\_guides.html#platform](https://bonseyes.gitlab.io/bonseyes-cli/pages/user_guides.html#platform)
- Developer guide: [https://bonseyes.gitlab.io/bonseyes-cli/pages/developer\\_guides.html#platform](https://bonseyes.gitlab.io/bonseyes-cli/pages/developer_guides.html#platform)

Developer platforms can be accessed upon request at <https://gitlab.com/bonseyes/platforms/>.

### 3.4. Benchmark framework

LPDNN provides a structured benchmarking layer to analyze the execution of AI Applications on HW platforms. LPDNN integrates a first analysis of static metrics during the offline compilation of the AI Applications. Metrics such as FLOPS, model parameters and model storage are obtained. Besides, LPDNN integrates two methods to benchmark the execution of AI Apps:

- **Built-in anchors:** C++ anchors are included in LPDNN to measure the latency and memory consumption of the AI Apps during their executions. Several anchors are included:
  - Pre-processing: This anchor measures the latency and memory of a given pre-processor:
    - Image: cropping, normalization, resize, filtering
    - Audio: mfcc feature computation
    - Signal: resizing, filtering
  - Inference: This anchor measures the latency and memory spent during the forward pass of the neural network for a given inference engine.
  - Post-processing: This anchor measures the latency and memory spent during the post-processing of the outputs of the forward pass, e.g., decoding, non-max suppression, bounding box forming.
  - Total execution: This anchor measures the latency and memory spent during the total execution of the AI App.
- **HW resources monitoring:** Python script that monitors the HW resources of the system while the execution of the AI App. Metrics such as CPU/GPU/NPU processor load, peak/average memory usage, memory bandwidth, temperature, etc.

Overall, these are the following metrics are LPDNN collects:

- **Static metrics (not measured on device):**
  - FLOPs
  - Parameters
  - Storage
- **Dynamic metrics (measured on device):**
  - Latency



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 101017142

- Throughput
- Average Memory (CPU/GPU)
- Peak Memory (CPU/GPU)
- Memory Bandwidth (CPU/GPU)
- Processor load (CPU/GPU/NPU)
- Power consumption
- Temperature

More details about the benchmark process will be given in D6.2.

## 4. Deployment of AI Apps and Benchmark

In this section, we showcase what elements are needed for deployment and how to execute a benchmark.

### 4.1. Deployment of AI Apps

To be able to execute AI applications on a hardware platform, two elements are required:

#### A. An LPDNN AI application:

It defines the class and structure of the AI application, the DNN models' architecture and weights, its deployment configuration and the pre- and post-processing that it takes. LPDNN AI applications are platform-specific, although the same AI application can be executed on different HW platforms as long as the selected inference engine or backends are supported by the HW platform.

#### B. An LPDNN Deployment Package:

Collection of tools, executables, libraries, inference engines and backends that allows the actual execution of the LPDNN AI application. The collection of libraries and binaries that are copied on the target platform for the execution of the DNN is called a *runtime*. The runtime dynamically loads an AI application and executes it based on its defined configuration. LPDNN's deployment packages are platform specific as they contain the inferences engines and backends supported by the HW platform.

#### 4.1.1 LPDNN AI Application (AI App)

An LPDNN AI App is composed of the following files:

- ❖ **ai\_app\_config.json**: This file is the main descriptor of an AI App. It defines the AI App's components and their type, e.g., image\_classification, object\_classification, face\_recognition, audio\_classification and signal\_processing, the type of pre- and post-processing as well as the inference engine to use to execute the DNN model. It also points to the model architecture and weights file.
- ❖ **ai\_app.yml**: This file defines the AI App metadata and license type. It also describes the platform, runtime and challenge that the AI App was initially compiled for. This file is not used by the runtime, but by other deployment tools.
- ❖ **DNN model**: A DNN model describes the model architecture and the trained weights. A DNN model may come on different forms based on the selected inference engine, e.g., model.json, model.param, model.bin or model.onnx.





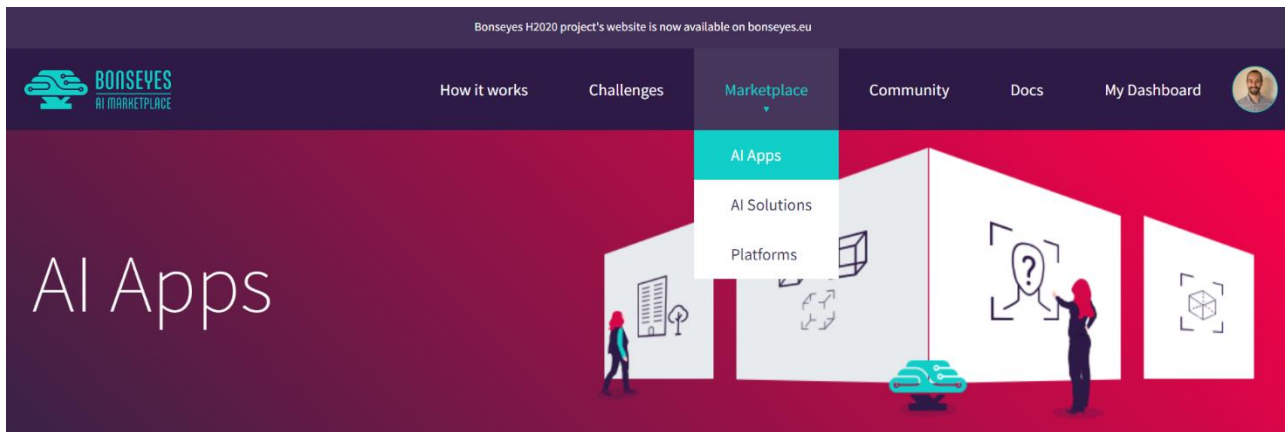


Figure 5 Bonseyes Marketplace Catalog

If the user would like to create its own application, detailed documentation is provided to generate a new one at [Bonseyes docs for AI-App generation](#).

#### 4.2.2. Requirements

To be able to execute AI applications on a hardware platform, the following steps need to have been performed:

1. Setup the local environment as explained in Bonseyes [Prerequisites](#) doc.
2. Set up the target hardware as explained in [Setup platform](#) section of Bonseyes' doc and have the `_${platformName}_src`, `_${platformName}_build` and `_${platformName}_config` folders in your machine.
3. Install python packages in the target board as explained in the [Packages's section](#).

Once those steps are completed, change directory to the folder where you built your target platform during the Setup platform section, e.g., my-bonseyes-platform.

#### 4.2.3. Download AI App and deployment package

New AI App can be generated by following the [Bonseyes docs for AI-App generation](#). Already built AI Apps can also be obtained from the Bonseyes Marketplace's catalogue. To download an AI-App from the Marketplace's catalogue execute the following command (replace `ai_app` by the name you would like to give to your AI-App):

```
bonseyes marketplace download-ai-app --output-dir ai_app
```

A dialog will prompts asking you to choose the AI-App you would like to download and the ai app will be downloaded in the directory `ai_app` as shown in Figure 6. The AI App will contain a series of files, as explained in Section 4.1.1.







AI-on-demand experimentation platform (Acumos) allows the creation of AI pipelines based on dockerized elements (Assets) that are interconnected by an orchestrator. Such containers may have any sort of artifact inside, e.g., dataset, AI model, or algorithm, and provide interfaces that define the API to communicate or transfer the content from one container to another. Users can create docker components and experiment with them in the form of an AI pipeline by defining an orchestrator to connect them up.

Partners of WP6 have held coordination meetings with AI4EU partners to understand the requirements and dependencies that are required to integrate the BenchmarkaaS into Acumos. Fig. 8 describes a potential interconnection between the Acumos and the benchmarkaaS:

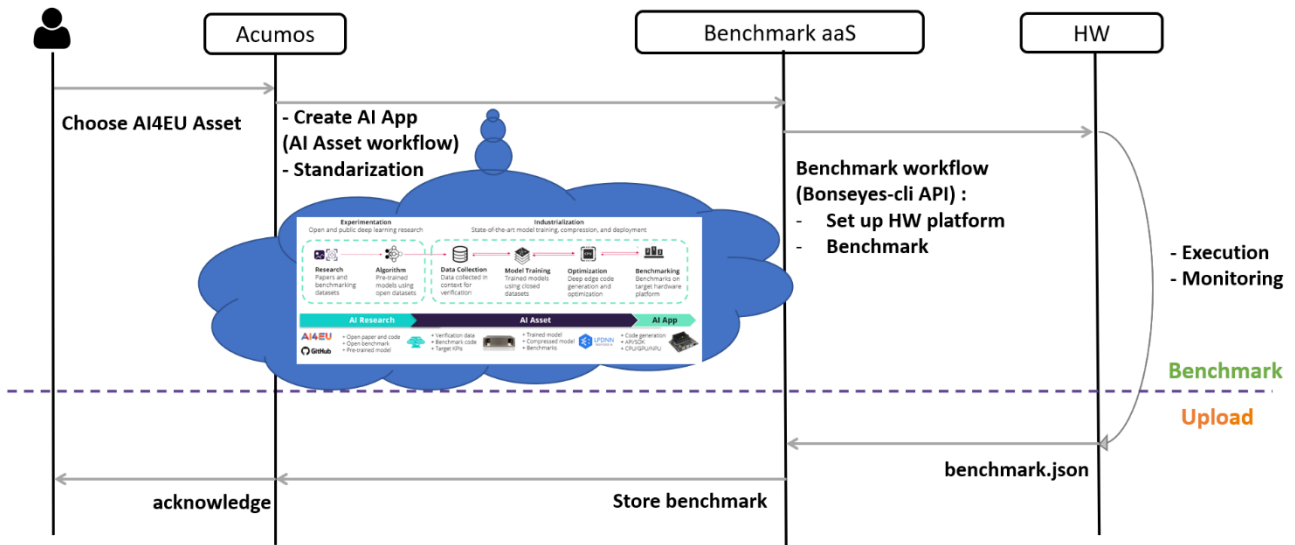


Figure 8 Benchmark as a Service

The understanding that has been reached so far indicates that a dockerized workflow (within the blue cloud in Fig.8) should be put in place to integrate the BenchmarkaaS into Acumos. The dockerized workflow that is proposed is the Bonseyes AI Asset workflow that should allow the following functions:

- Align docker dependencies and interfaces between Acumos and LPDNN’s docker environments.
- Allow the transfer of AI models from Acumos’s containers (Assets) to the Bonseyes’s ecosystem by providing interfaces that allow Acumos’s Assets to be “benchmarkable” by the BenchmarkaaS.
- Transform the AI models into LPDNN AI Applications that can be executed by the benchmark layer.

More details and efforts will be put into the integration/compatibility of the benchmark as a Service with the AI-on-demand platform during the second half of the StairwAI project, as part of T2.5. However, to go beyond the initial proposal in the Grant agreement, appropriate resources should be mobilised to this effect.





## 5. Conclusion and future work

This document has presented the first version of the *Benchmark software framework* proposed within WP6 Task 6.1 of StairwAI project. The document describes how LPDNN inference framework has been used and extended to accommodate a benchmark layer that allows to profile AI Applications on heterogeneous HW platforms. The deployment and benchmark of AI Applications has been showcased, illustrating the flow and commands that a user need to execute to benchmark an AI Application.

As future work for the 2<sup>nd</sup> release of the *Benchmark software framework*, more efforts will be made to integrate Huawei's embedded platform and inference engine into LPDNN, increasing the capabilities for profiling AI Algorithms on a wide range of platforms,

In addition, more effort will be put into the integration/interoperability of the benchmark as a Service with the AI4EU platform, providing the Bonseyes AI Asset workflows as a bridge to allow the benchmark of a variety of AI4EU's Assets.

## Bibliography

1. Prado, Miguel De, et al. "Bonseyes AI pipeline—Bringing AI to you: End-to-end integration of data, algorithms, and deployment tools." *ACM Transactions on Internet of Things* 1.4 (2020): 1-25.

