



## Design and Runtime Framework for Accelerating the Development of AI Applications in the Computing Continuum

Francesco Lattari, Politecnico di Milano

[francesco.lattari@polimi.it](mailto:francesco.lattari@polimi.it)



AI-SPRINT project has received funding from the European Union Horizon  
2020 research and innovation programme under Grant Agreement **No. 101016577**.

COORDINATOR



**POLITECNICO**  
MILANO 1863



**BSC** Barcelona  
Supercomputing  
Center  
Centro Nacional de Supercomputación



**TECHNISCHE  
UNIVERSITÄT  
DRESDEN**



**UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA**



**GRENOIRE**



**Beck  
et al.** work.  
together.



**CLOUD  
& HEAT**



**7** bulls.com



**TAnalysis**

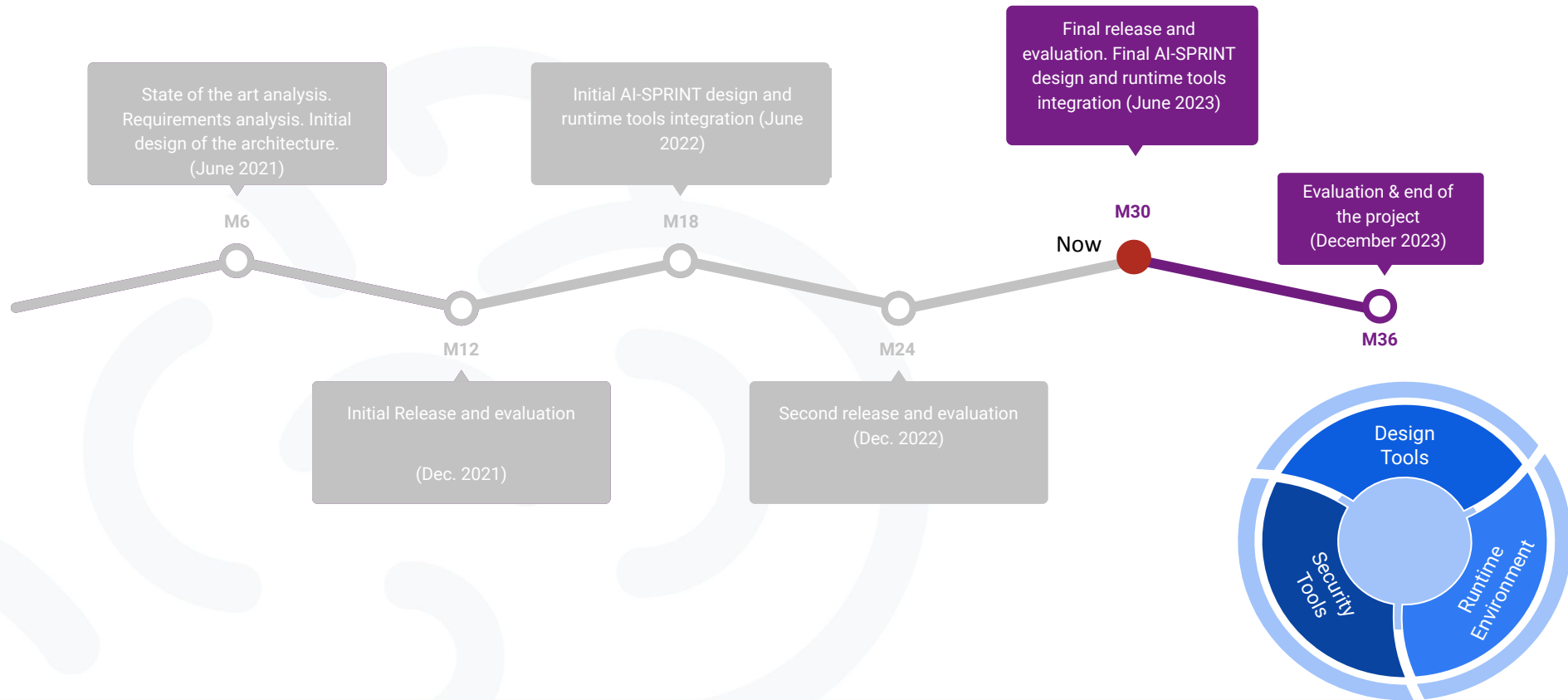


**Trust-IT Services**  
communicating to markets

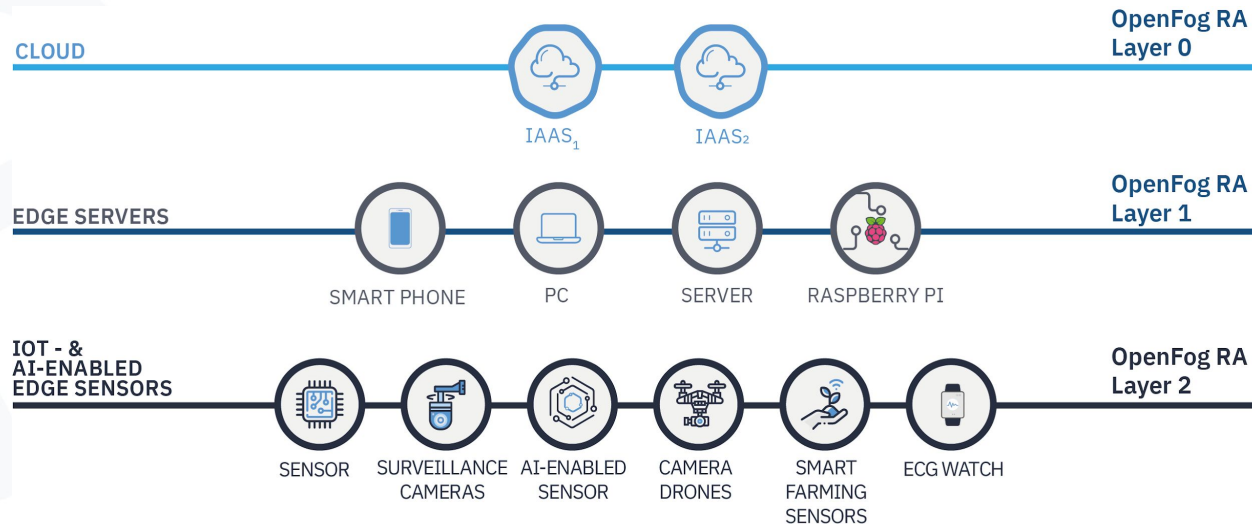


**IDC**

# AI-SPRINT Development Roadmap



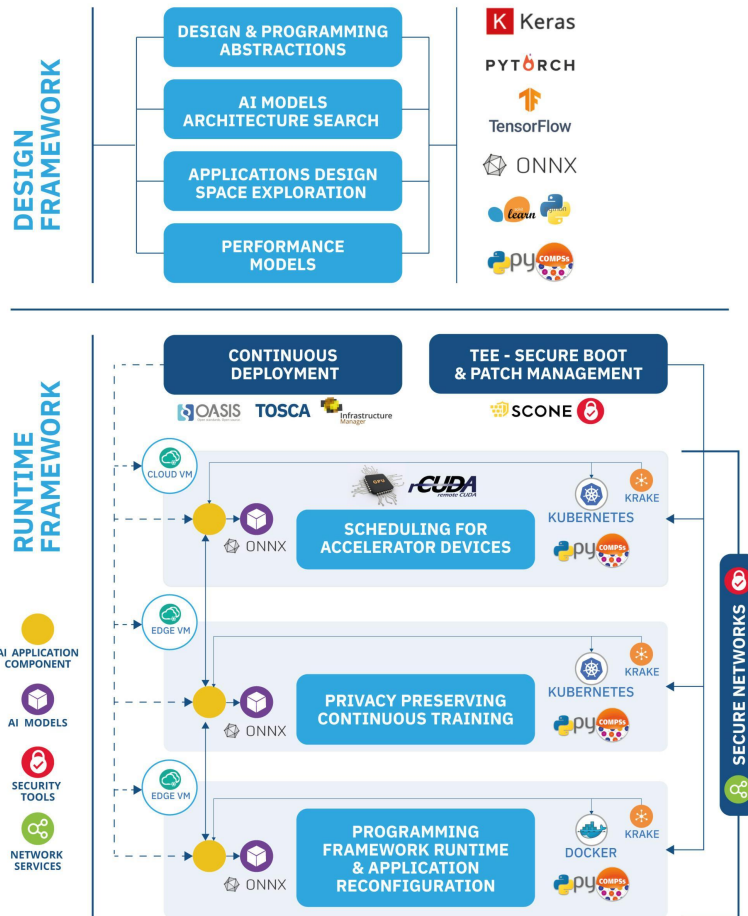
- By 2026, AI worldwide market will approach \$900 billion (CAGR 18.6%<sup>1</sup>) while edge computing will reach \$324 billion (CAGR 13.6%<sup>2</sup>)
- AI needs resources at the edge of the network
- New challenges from the infrastructural perspective



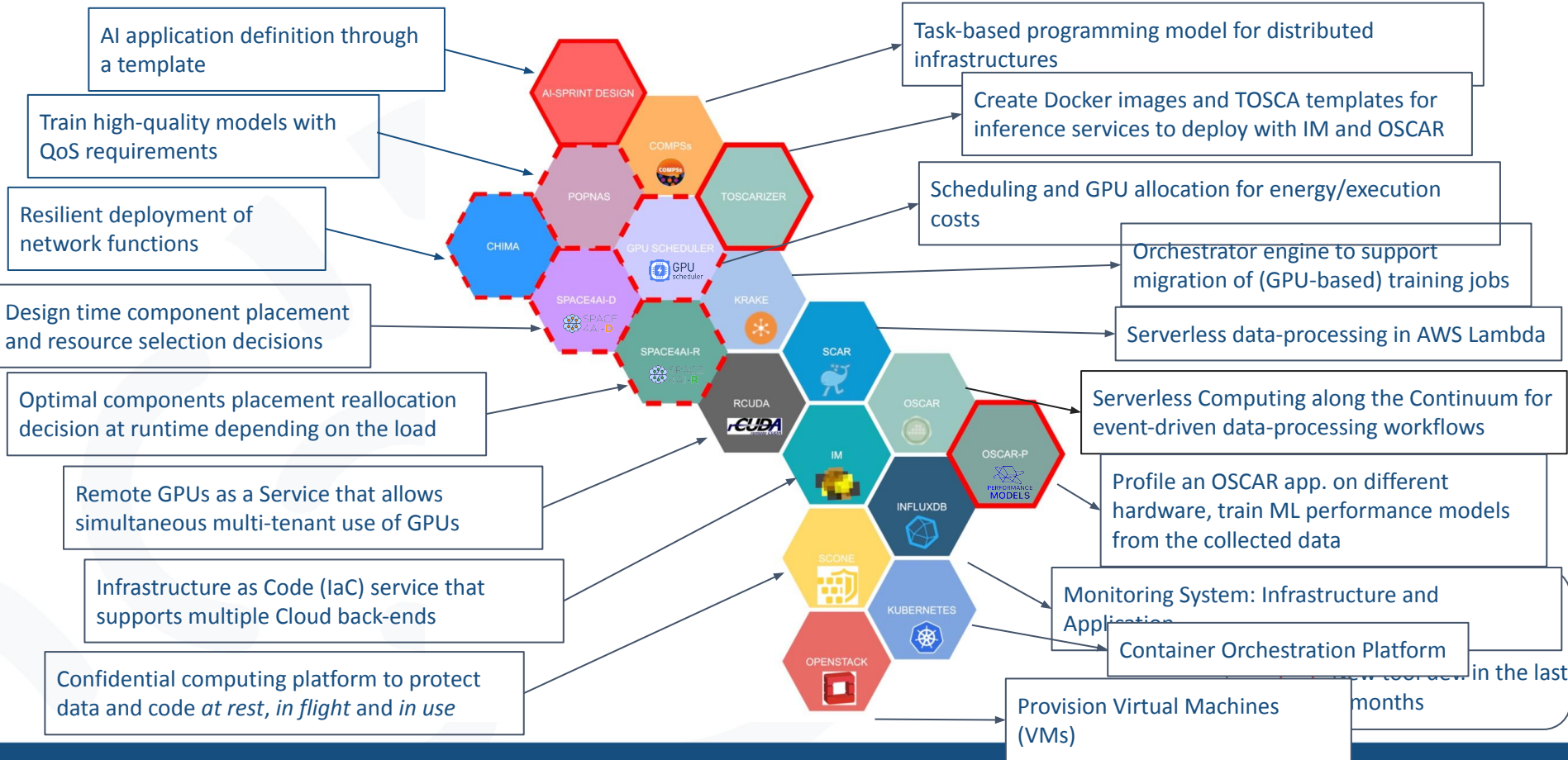
<sup>1</sup>IDC Semiannual Artificial Intelligence Tracker, July 2022

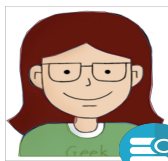
<sup>2</sup>IDC Worldwide Edge Spending Guide, August 2022

- Simplified programming models
- Automated deployment and dynamic reconfiguration
- Secure execution of AI applications
- Highly specialized building blocks for privacy preservation, distributed training, and architecture enhancement
- Open source



# Main Technological Components





**Application Developer**



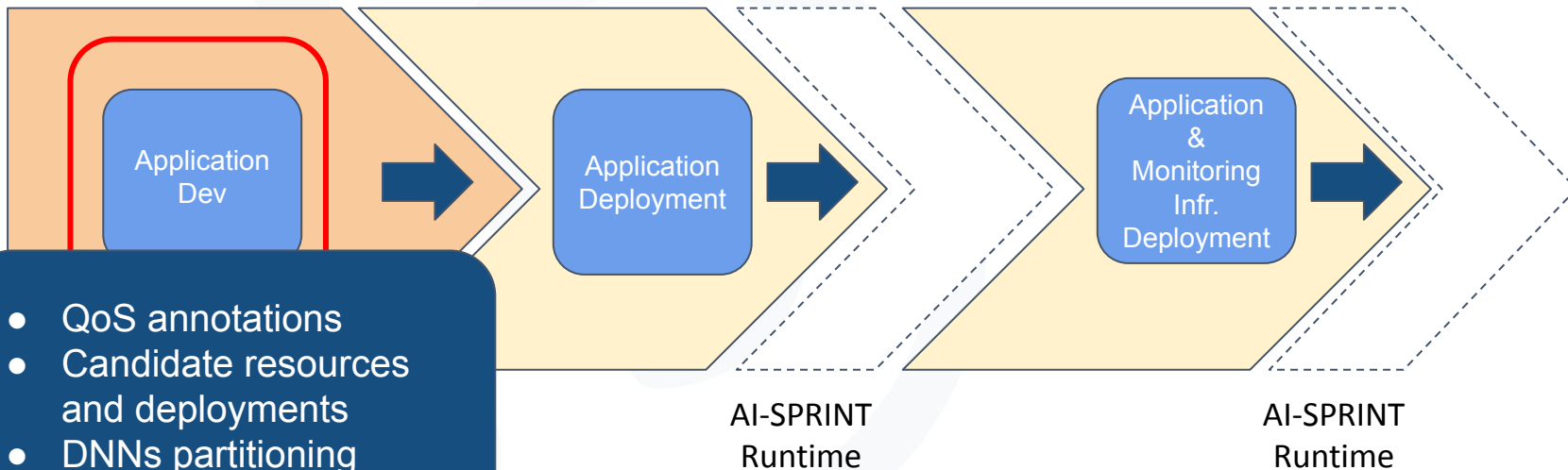
**Application Architect**



**Application Manager**

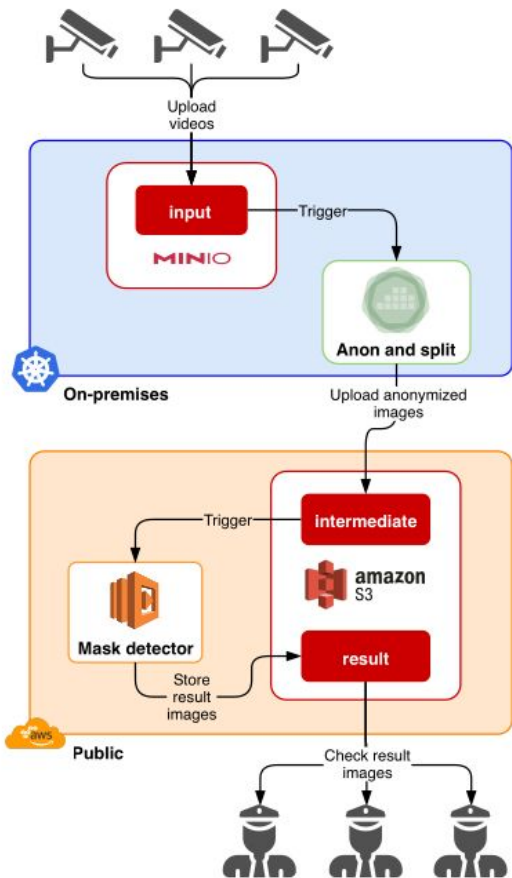


**Application Manager**

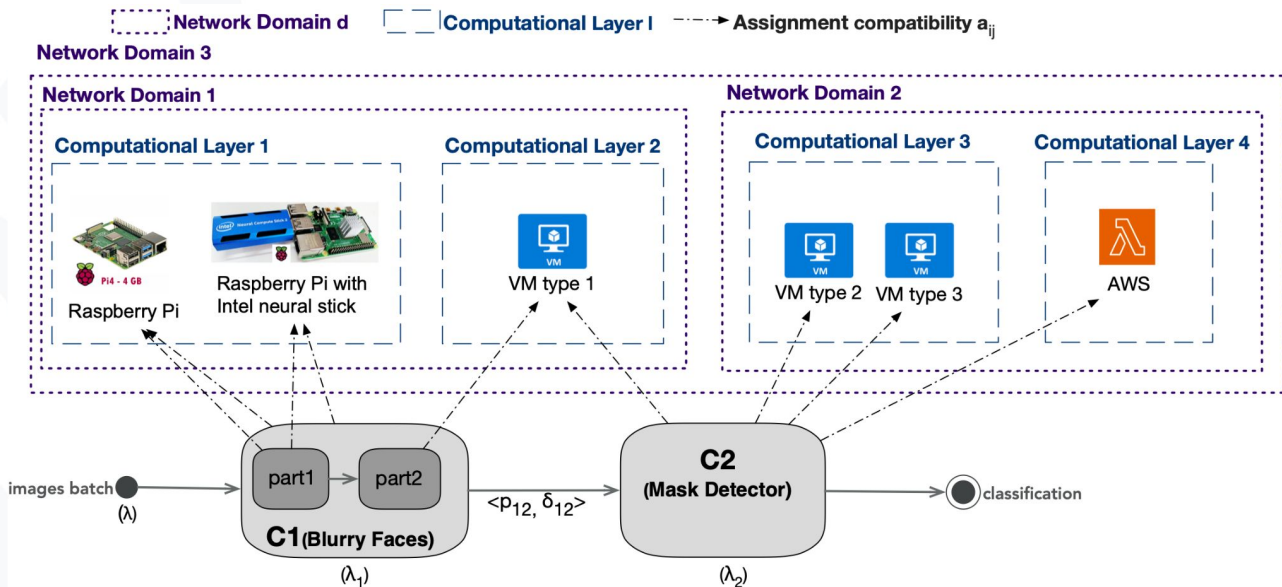


- QoS annotations
- Candidate resources and deployments
- DNNs partitioning

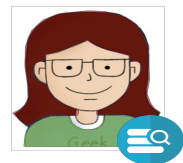




- Single or multiple Python applications, e.g., “Anon and split” and “Mask detector”
- Single or multiple candidate resources for each component
- Workflow execution orchestrated by **OSCAR** (<https://oscar.grycap.net>)

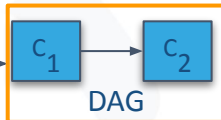






**Application Developer**

Annotated Python code



**Application Manager**

Candidate Resources



**Application Architect**

Candidate Deployments

```

AI-SPRINT_MASK_DETECTION_APP
├── aisprint
│   ├── deployments
│   ├── designs
│   ├── logs
│   ├── ams
│   └── common_config
│       ├── application_dag.yaml
│       ├── candidate_deployments.yaml
│       └── candidate_resources.yaml
├── im
├── oscar
├── pycompss
├── space4ai-d
├── src
├── blurry-faces-onnx
│   ├── onnx
│   │   └── version-RFB-640.onnx
│   ├── main.py
│   ├── requirements.txt
│   └── utils.py
├── mask-detector
│   ├── cfg
│   │   ├── obj.names
│   │   ├── yolov3-tiny_obj_test.cfg
│   │   ├── yolov3-tiny_obj_train_tiny8.weights
│   │   ├── yolov3-tiny_obj_train.cfg
│   │   ├── main.py
│   │   └── requirements.txt

```

```

@Component(name='mask-detector')
@exec_time(local_time_thr=10)
@device_constraints(ram=1024, vram=2048)
@security(trustedExecution=False, networkShield=False, filesystemShield=False)
def main(args):

```

```

System:
  name: ai-sprint_mask_detection_app
  components: ['blurry-faces-onnx', 'mask-detector']
  dependencies: [['blurry-faces-onnx', 'mask-detector', 1]]

```

```

System:
  name: Mask Detection Application
  NetworkDomains:
    ND1:
      name: Network Domain 1
      AccessDelay: 0.00000277
      Bandwidth: 40000
      subNetworkDomains: []
      ComputationalLayers:
        computationalLayer1:
          name: Edge Layer
          number: 1
          type: PhysicalAlreadyProvisioned
      Resources:
        resource1:
          name: RaspPi
          totalNodes: 3
          description: Raspberry PI
          cost: 0.6
          memorySize: 4096
          operatingSystemDistribution:
          operatingSystemType: Linux
          operatingSystemVersion: 10
          operatingSystemImageId: N

```

```

Components:
  component1:
    name: blurry-faces-onnx
    candidateExecutionLayers: [1,2]
    Containers:
      container1:
        image: registry.gitlab.polimi.it/ai-sprint/blurry-faces-onnx
        memorySize: 2048
        computingUnits: 0.9
        trustedExecution: False
        networkProtection: False
        filesystemProtection: False
        GPURequirement: False
        candidateExecutionResources: [RaspPi]
      container2:
        image: registry.gitlab.polimi.it/ai-sprint/mask-detector
        memorySize: 2048
        computingUnits: 0.9
        trustedExecution: False
        networkProtection: False
        filesystemProtection: False
        GPURequirement: False
        candidateExecutionResources: [VM1]

```

Annotation name
<i>@aisprint.component_name</i>
<i>@aisprint.exec_time</i>
<i>@aisprint.expected_throughput</i>
<i>@aisprint.partitionable_model</i>
<i>@aisprint.device_constraints</i>
<i>@aisprint.early_exits_model</i>
<i>@aisprint.model_performance</i>
<i>@aisprint.detect_metric_drift</i>
<i>@aisprint.security</i>

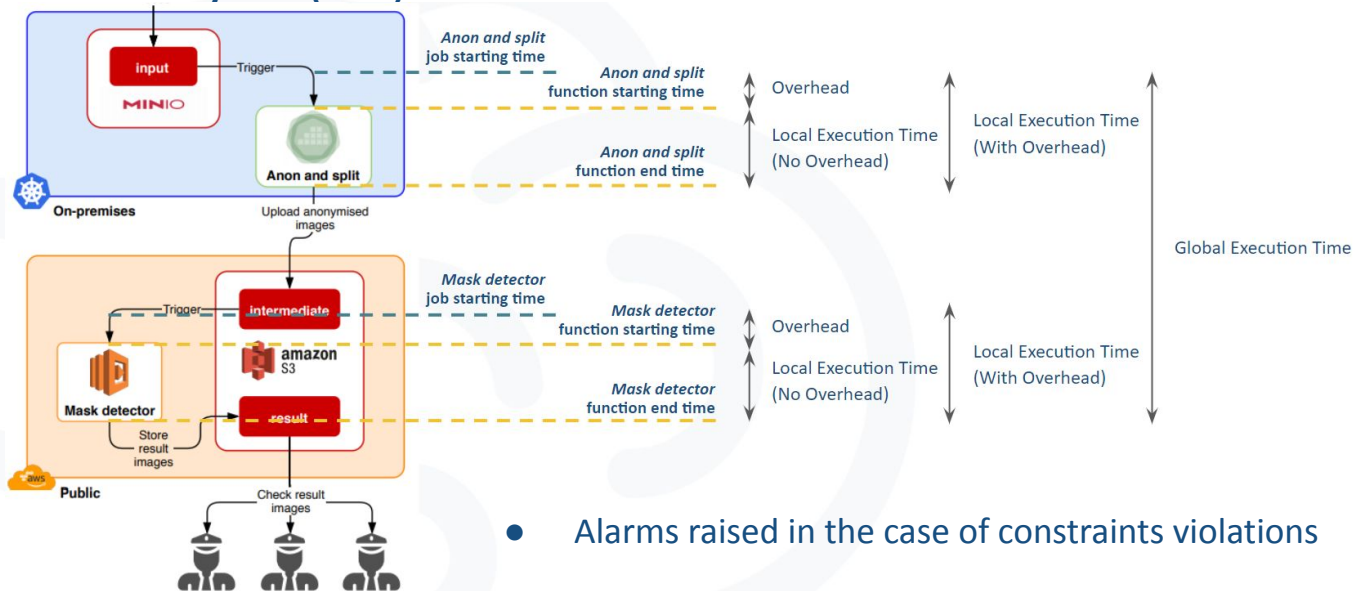
Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
@aisprint.partitionable_model
@aisprint.device_constraints
@aisprint.early_exits_model
@aisprint.model_performance
@aisprint.detect_metric_drift
@aisprint.security

```
def exec_time(local_time_thr, global_time_thr, prev_components)
```

Allows users defining execution time constraints for single (local) or multiple components (global)

- Constraints are automatically monitored by the **AI-SPRINT Monitoring Subsystem (AMS)**

<https://github.com/aimlhubio/ai-sprint-monitoring>



- Alarms raised in the case of constraints violations

Annotation name
<i>@aisprint.component_name</i>
<i>@aisprint.exec_time</i>
<b><i>@aisprint.expected_throughput</i></b>
<i>@aisprint.partitionable_model</i>
<i>@aisprint.device_constraints</i>
<i>@aisprint.early_exits_model</i>
<i>@aisprint.model_performance</i>
<i>@aisprint.detect_metric_drift</i>
<i>@aisprint.security</i>

```
def expected_throughput(rate)
```

Allows the users to define the expected application throughput

- i.e., the expected invocation rate expressed as the number of invocations per time unit (number of invocations per second)

```
def partitionable_model(onnx_file, num_partitions)
```

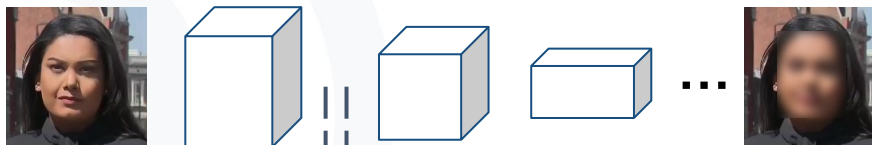
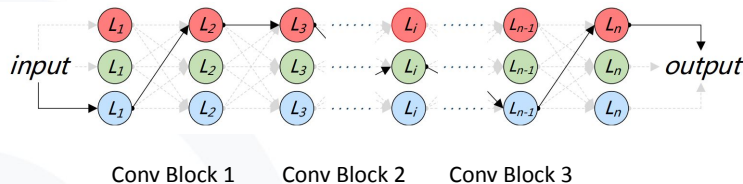
Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
<b>@aisprint.partitionable_model</b>
@aisprint.device_constraints
@aisprint.early_exits_model
@aisprint.model_performance
@aisprint.detect_metric_drift
@aisprint.security

Gives the user the possibility of defining a *partitionable* Deep Neural Network (DNN), provided using the **Open Neural Network Exchange (ONNX)** format

- i.e., the model is divided in different parts, which can be executed on different computational layers (L)



<https://onnxruntime.ai/>



- Automatic split performed by the SPACE4AI-D-Partitioner tool (<https://gitlab.polimi.it/ai-sprint/ai-sprint-design/-/tree/master/src/aisprint/space4aidpartitioner>)

Annotation name
<i>@aisprint.component_name</i>
<i>@aisprint.exec_time</i>
<i>@aisprint.expected_throughput</i>
<i>@aisprint.partitionable_model</i>
<b><i>@aisprint.device_constraints</i></b>
<i>@aisprint.early_exits_model</i>
<i>@aisprint.model_performance</i>
<i>@aisprint.detect_metric_drift</i>
<i>@aisprint.security</i>

```
def device_constraints(ram, vram)
```

The annotation allows the users to specify a set of minimal resources that the device on which the component will be deployed must have. In particular the

- **ram**: minimum memory in GB required to run the annotated component
- **vram**: minimum video memory in GB required to run the annotated component. A  $vram > 0$  implicitly highlights the need for a Graphics Processing Unit (GPU).



```
def early_exits_model(onnx_file, condition_function, transition_probabilities)
```

Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
@aisprint.partitionable_model
@aisprint.device_constraints
<b>@aisprint.early_exits_model</b>
@aisprint.model_performance
@aisprint.detect_metric_drift
@aisprint.security

Similarly to the partitionable models, allows automatic partitioning of DNN with early exits

- Execution can be stopped earlier in the network based on user-defined conditions. E.g., [BranchyNet: Fast Inference via Early Exiting from Deep Neural Networks](#)

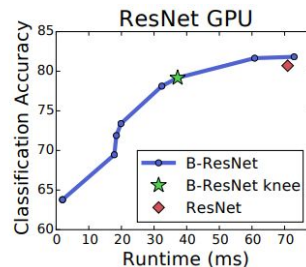
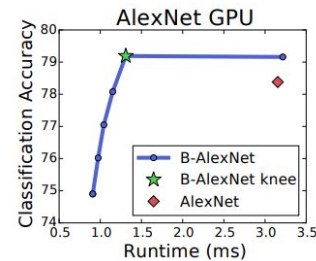
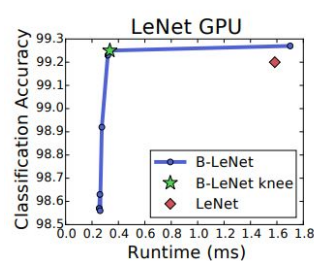
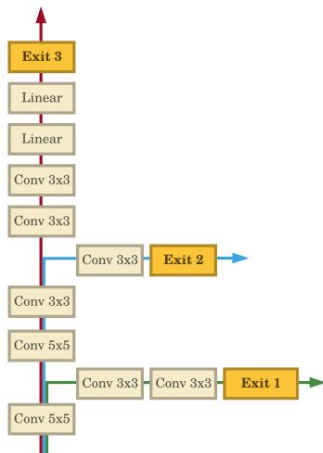


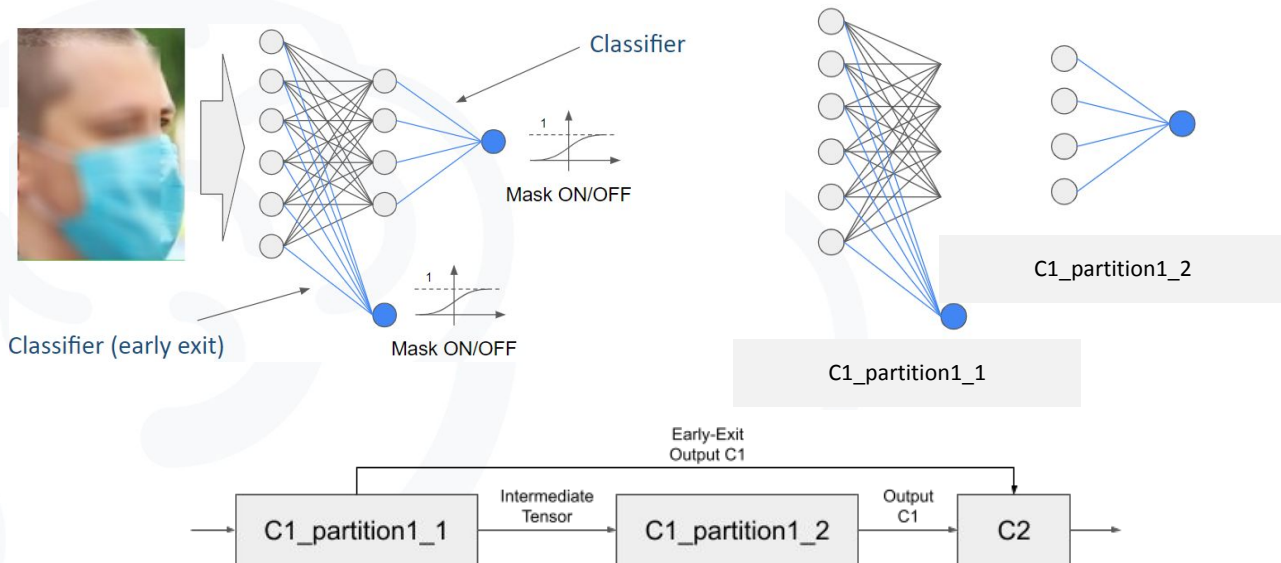
Fig. 1: A simple BranchyNet with two branches added to the baseline (original) AlexNet. The first branch has two convolutional layers and the second branch has 1 convolutional layer. The “Exit” boxes denote the various exit points of BranchyNet.

```
def early_exits_model(onnx_file, condition_function, transition_probabilities)
```

Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
@aisprint.partitionable_model
@aisprint.device_constraints
<b>@aisprint.early_exits_model</b>
@aisprint.model_performance
@aisprint.detect_metric_drift
@aisprint.security

Similarly to the partitionable models, allows automatic partitioning of DNN with early exits

- AI-SPRINT automatically split the DNN-based components at early exits, generating new components corresponding to the network segments

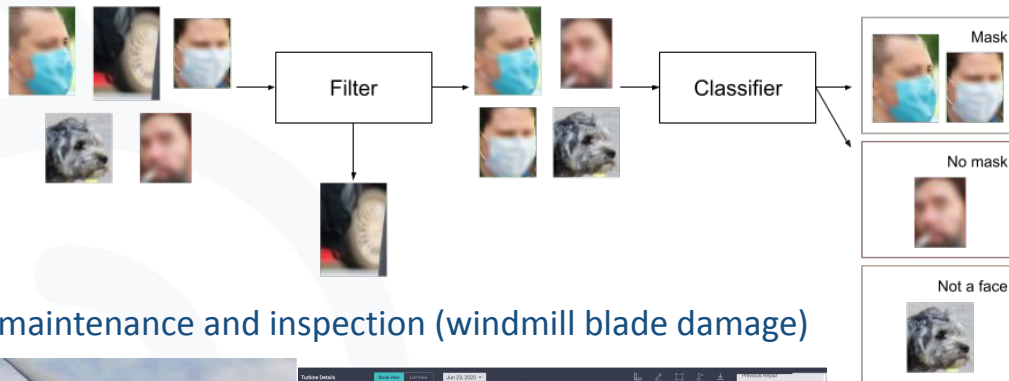


```
def model_performance(metric, metric_thr, filtered_class)
```

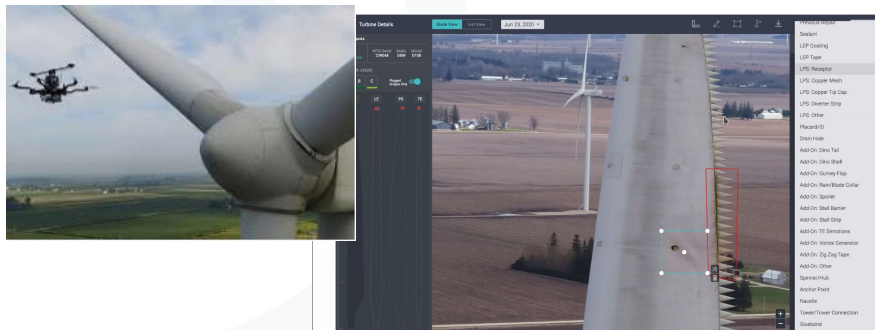
Annotation name
<i>@aisprint.component_name</i>
<i>@aisprint.exec_time</i>
<i>@aisprint.expected_throughput</i>
<i>@aisprint.partitionable_model</i>
<i>@aisprint.device_constraints</i>
<i>@aisprint.early_exits_model</i>
<i>@aisprint.model_performance</i>
<i>@aisprint.detect_metric_drift</i>
<i>@aisprint.security</i>

Allows defining applications with degraded performance

- Use case: filter+classifier applications



- E.g., maintenance and inspection (windmill blade damage)



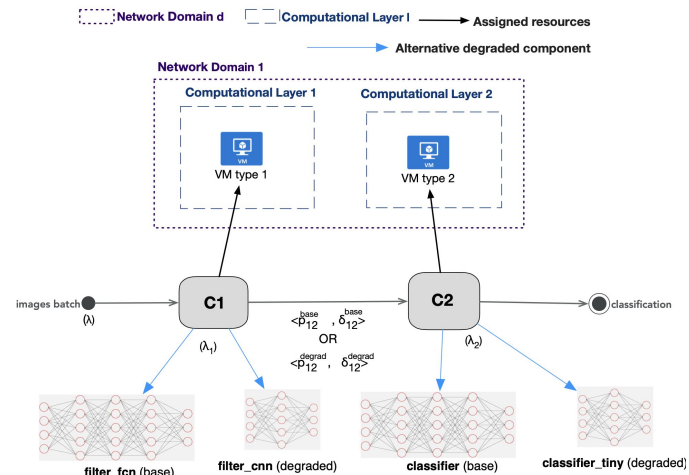
```
def model_performance(metric, metric_thr, filtered_class)
```

## Allows defining applications with degraded performance

- The user can provide multiple version of both filter and classifier components. Performance of the alternative applications is automatically computed and applications are ordered accordingly.
- SPACE4AI-R Runtime tool is able to switch from one alternative workflow to another in the case is needed

Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
@aisprint.partitionable_model
@aisprint.device_constraints
@aisprint.early_exits_model
@aisprint.model_performance
@aisprint.detect_metric_drift
@aisprint.security

```
! alternative_workflows.yaml x
common_config > ! alternative_workflows.yaml > {} System > [ ] alternative_depend
1 System:
2   name: filter_classifier_app
3   alternative_dependencies:
4     - alternative_1:
5       dependency: [['filter', 'classifier', 1]]
6       metric:
7         name: average_f1
8         value: 0.9
9     - alternative_2:
10      dependency: [['filter_v2', 'classifier_v2', 1]]
11      metric:
12        name: average_f1
13        value: 0.85
14
15     - alternative_6:
16      dependency: [['filter_v3', 'classifier_v3', 1]]
17      metric:
18        name: average_f1
19        value: 0.6
```

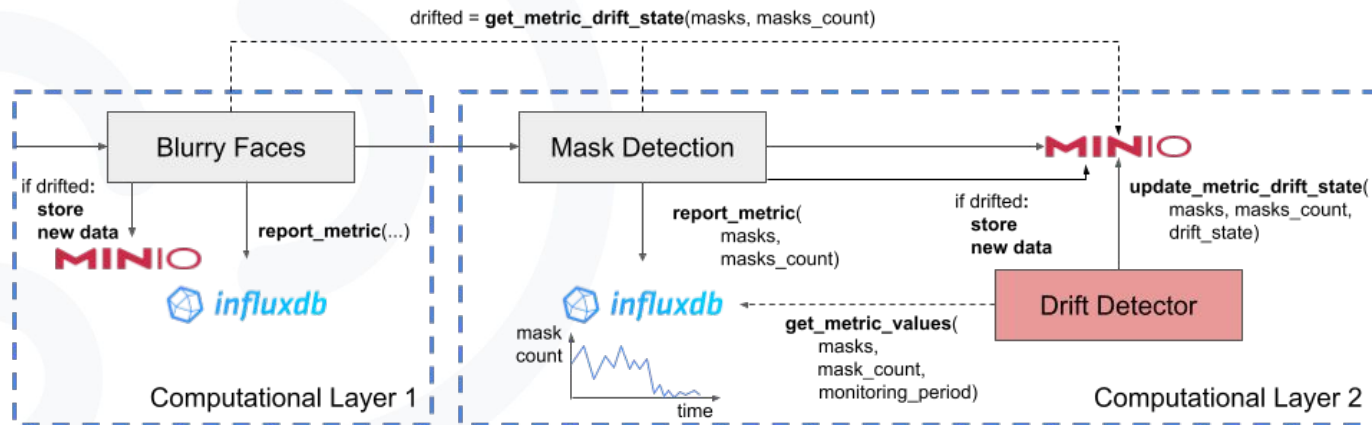


Annotation name
@aisprint.component_name
@aisprint.exec_time
@aisprint.expected_throughput
@aisprint.partitionable_model
@aisprint.device_constraints
@aisprint.early_exits_model
@aisprint.model_performance
<b>@aisprint.detect_metric_drift</b>
@aisprint.security

```
def detect_metric_drift(metric, field, statistical_test,
                       test_threshold, detection_interval,
                       monitoring_period, data_collection_period)
```

Allows detecting data drift at runtime, by triggering the automatic deployment of the **Drift Detector** tool

- Periodically queries user-defined metrics (stored in InfluxDB) and run statistical algorithms to detect changes in the time series
- Allows collecting new data after the drift to be used for re-training the DNN-based components



Annotation name
<code>@aisprint.component_name</code>
<code>@aisprint.exec_time</code>
<code>@aisprint.expected_throughput</code>
<code>@aisprint.partitionable_model</code>
<code>@aisprint.device_constraints</code>
<code>@aisprint.early_exits_model</code>
<code>@aisprint.model_performance</code>
<code>@aisprint.detect_metric_drift</code>
<b><code>@aisprint.security</code></b>

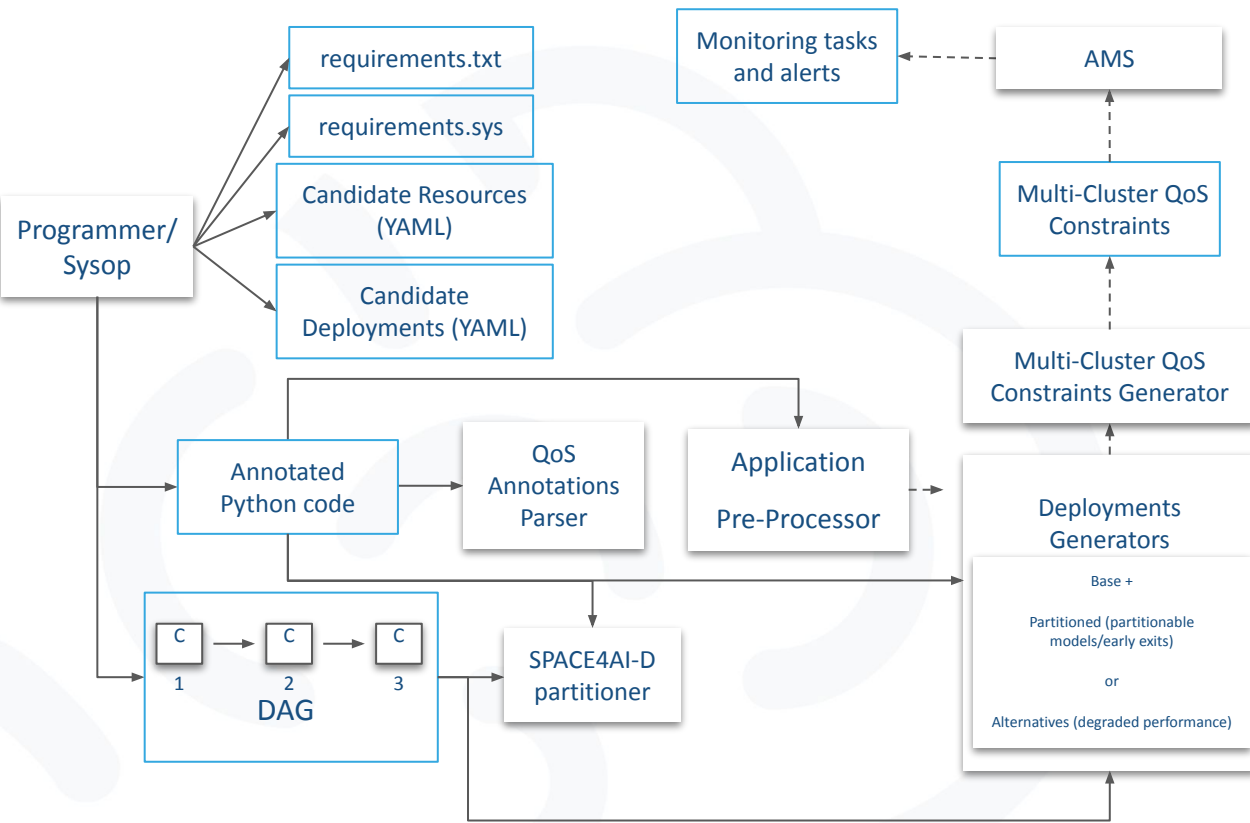
```
def security(trustedExecution, networkShield, filesystemShield,  
            confProc, integrityProc, confRest)
```

Allows users annotating components to receive security guarantees while executing:

- Enables *trusted execution environments* and secure boot
- Wraps TCP connections using the **SCONE** network shielding layer
- Enables the encryption of all the files written and read by the process executing the task



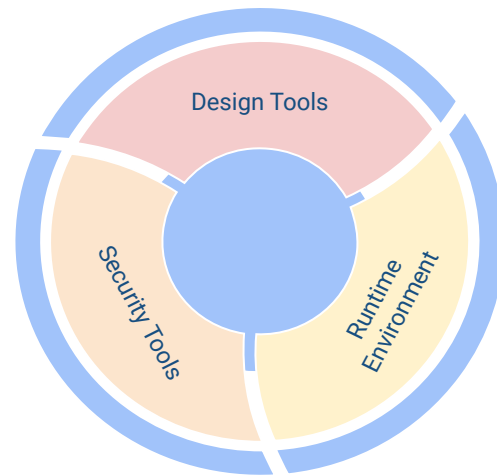
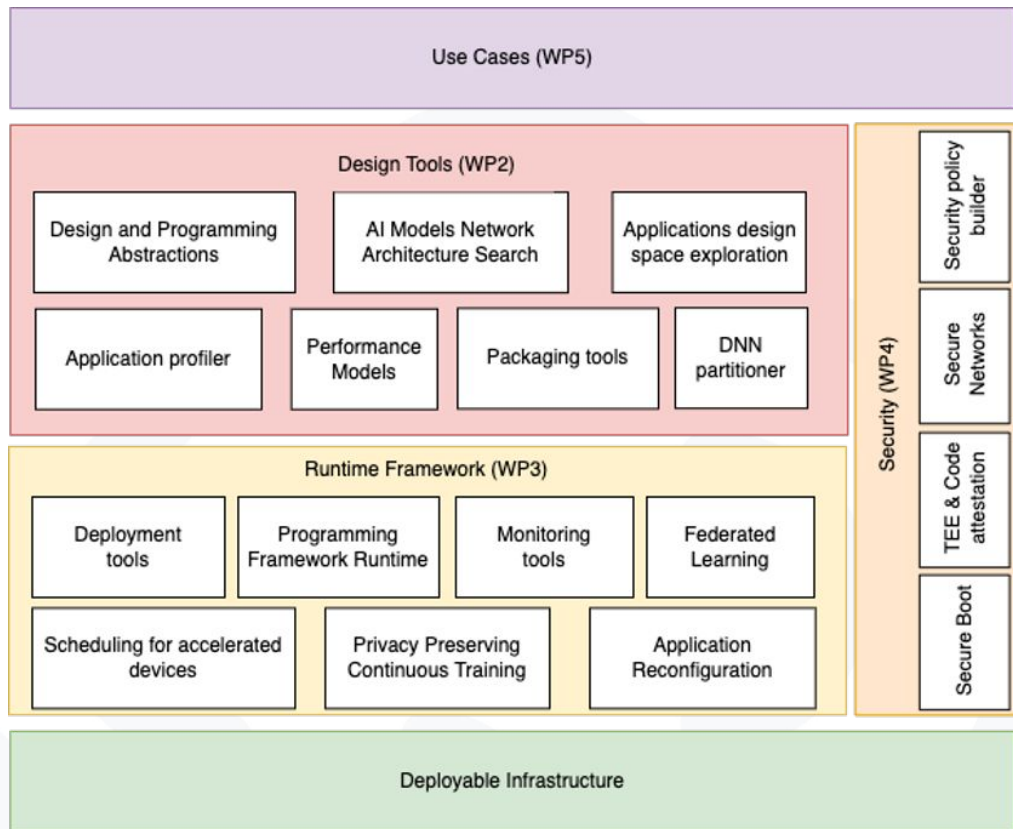
Run AI-SPRINT Design with AI-SPRINT Studio <https://gitlab.polimi.it/ai-sprint/ai-sprint-studio>

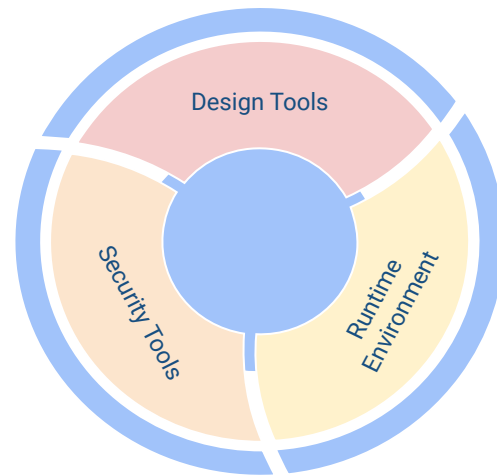
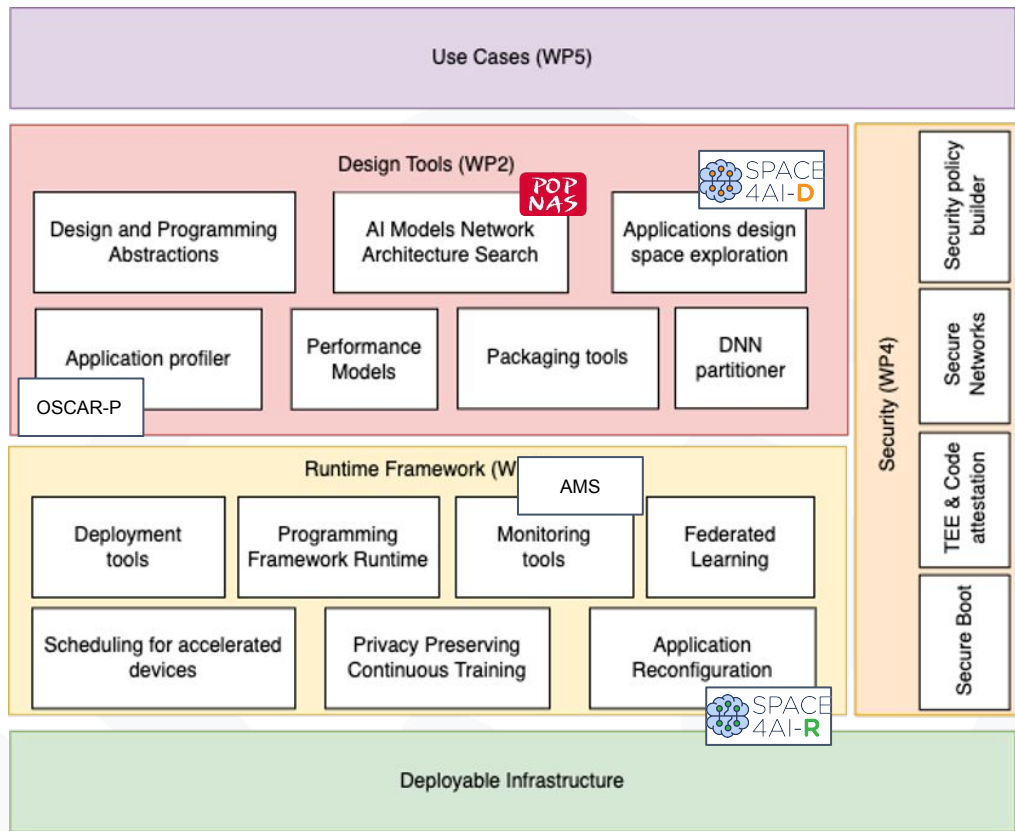


```

deployment2:
  ExecutionLayers:
    1:
      components:
        - blurry-faces-onnx
      local_constraints: {}
      global_constraints: {}
      throughput_component: blurry-faces-onnx
    2:
      components:
        - mask-detector-onnx_partition1
      local_constraints: {}
      global_constraints: {}
      throughput_component: blurry-faces-onnx
    3:
      components:
        - mask-detector-onnx_partition2
      local_constraints: {}
      global_constraints:
        global_constraint 1:
          path_components:
            - blurry-faces-onnx
            - mask-detector-onnx_partition1
            - mask-detector-onnx_partition2
          threshold: 30
      throughput_component: blurry-faces-onnx
deployment3:
  ExecutionLayers:
    1:
      components:
        - blurry-faces-onnx_partition1
      local_constraints: {}
      global_constraints: {}
      throughput_component: blurry-faces-onnx_partition1
    2:
      components:
        - blurry-faces-onnx_partition2
        - mask-detector-onnx
      local_constraints: {}
      global_constraints:
        global_constraint 1:
          path_components:
            - blurry-faces-onnx_partition1
            - blurry-faces-onnx_partition2
            - mask-detector-onnx
          threshold: 30
      throughput_component: blurry-faces-onnx_partition1
  
```









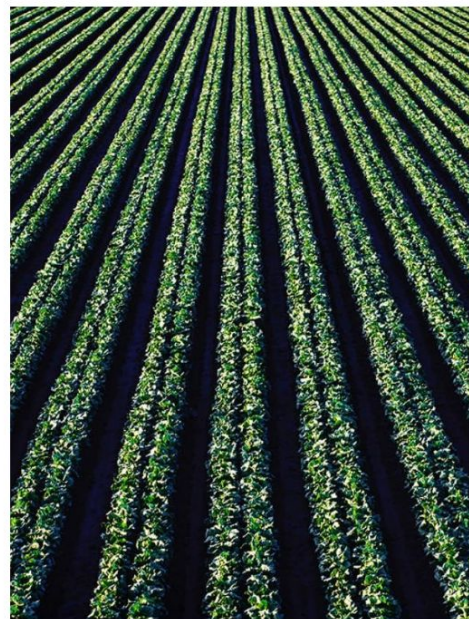
## Personalised Healthcare

Developing an automated system for personalised stroke risk assessment and prevention.



## Maintenance & Inspection

Creating an infrastructure that reduces downtime and revenue losses caused by degenerative asset performance.



## Farming 4.0

Delivering edge and intelligent sensors to optimise phytosanitary treatments.



The AI-SPRINT Alliance is composed of a group of a specialised supply and demand community of Software houses, AI-application developers, System integrators, Cloud Providers, Digital Innovation Hubs, and R&D initiatives that can use the **AI-SPRINT components and tools with technical support** from the project's partners





## COORDINATOR



POLITECNICO  
MILANO 1863



Barcelona  
Supercomputing  
Center  
Centro Nacional de Supercomputación



TECHNISCHE  
UNIVERSITÄT  
DRESDEN



UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA

GRENOBLE

Beck  
et al. *work.  
together.*

CLOUD  
& HEAT

7 bulls.com

TTAnalysis

Trust-IT Services  
*communicating to markets*

IDC

Thanks for your  
attention!

AI-SPRINT demos available at: <https://gitlab.polimi.it/ai-sprint/ai-sprint-examples>  
AI-SPRINT Studio library available at: <https://gitlab.polimi.it/ai-sprint/ai-sprint-studio>  
AI-SPRINT Studio Docker image available at: [registry.gitlab.polimi.it/ai-sprint/ai-sprint-studio](https://registry.gitlab.polimi.it/ai-sprint/ai-sprint-studio)