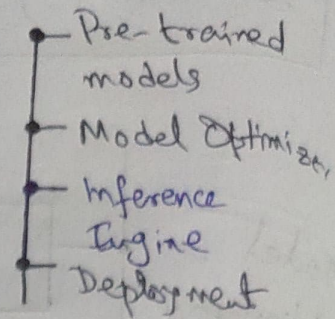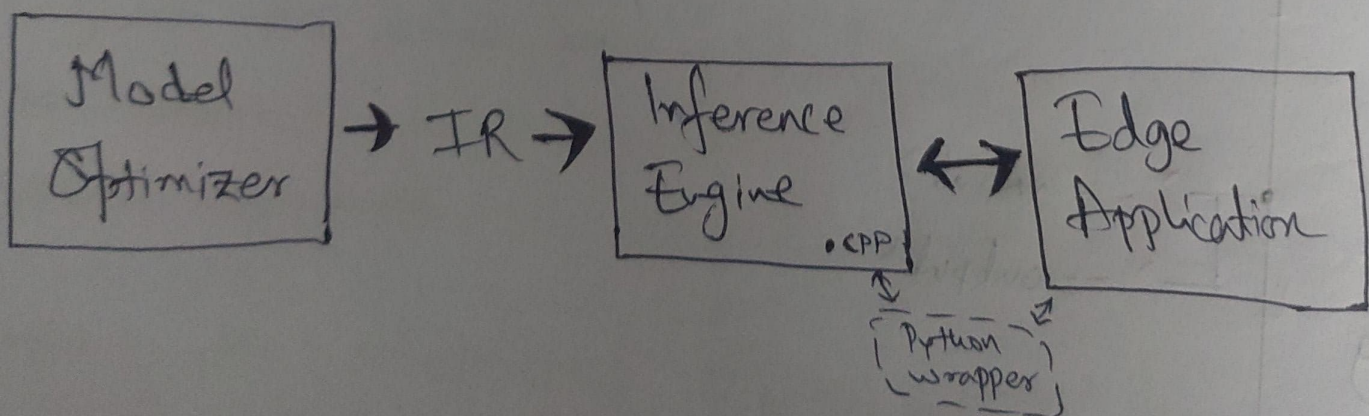# ④ THE INFERENCE ENGINE

- The Inference Engine runs the actual inference on the model.

- It only work w/ models that are:
  - already in the IR format as a pre-trained model, ~~model~~ OR
  - models converted to the IR format from the model optimizer.

The Inference Engine provides a library of CV functions needed for ~~the~~ inference.

```
┌──────────┐      ┌──────────┐        ┌──────────┐
│ Model    │ →IR→ │ Inference│  ↔     │ Edge     │
│ Optimizer│      │ Engine   │        │Application│
└──────────┘      │     •CPP │        └──────────┘
                  └────┬─────┘
                       ↕
                  ┌─────────┐
                  │ Python  │
                  │ wrapper │
                  └─────────┘
```

| Devices Supported by the Inference Engine |
| --- |

① CPUs &

② GPUs

③ FPGAs or Field Programmable Gate Arrays

④ VPUs or Visual Processing Units

    Eg: Intel Neural Compute Stick (NCS)

✳ All these need to be Intel hardware only.

| Using the Inference Engine with an IR |
| --- |

- Library ┬ openvino.inference_engine
     ├ class — IECore
     └ class — IENetwork

- IECore is the Python wrapper to work w/ the Inference Engine.

- IENetwork will initially hold the network & get loaded into IECore.

# Steps to load an IR into an IE

① Get the IR model (.xml) and weights (.bin), and pass them to the program.

② In the program, import the openvino IE library and from it, the classes IE Core and IENetwork.

③ Create an IE Core instance (say iecore)

④ Create an IENetwork instance (say net)

⑤ Add a CPU extension if needed

⑥ Get the supported layers of the network

⑦ from this list, check and see if there is any unsupported layer by comparing with all layers from net.layers.keys(). If any gets found, perform proper handling of this error

⑧ Finally, load the ie core network by providing the net and device name

Methods used in each step (refer docs)

• ③ IE Core()          ⑥ iecore.query_network(...)
  ④ IENetwork(...)      ⑦ net.layers.keys(), exit(...)
  ⑤ iecore.add_extension(...)  ⑧ iecore.load_network(...)

## Sending Inference Request to the IE

- The load-network method of IECore class returns an ExecutableNetwork object.

- The inference requests are made to this object.

- Types of requests:

| Synchronous | Asynchronous |
|---|---|
| (i) The app sits and waits for the infer-ence. | The The app can perform other tasks while making the inference. |
| (ii) Method used: infer() | Methods used: start_async() wait() |
| (iii) The main thread is "blocked". | Does not block any thread as the response may be slow. |
| (iv) The next frame is not gathered until the current frame's request is complete. | Sends a frame for inference while processing goes on in the next frame waiting f inference result. |

# Handling Results & App Integration

- ⊗ Attributes of InferRequest instances:

  (i) inputs → the image frame

  (ii) outputs → the results

  (iii) latency → the inference time of current request

- All inference requests are stored in a requests attribute in ExecutableNetwork

- To fetch an o/p:

  exec_net.requests [request_id]. outputs [output_blob]

- For app integration, we would need all the concepts learned so far, and more.

## Summary

**Inference Engine** — performs inference on models in the IR format

**Supported Devices**
- Intel proprietary hardware
- Eg. CPUs, GPUs, FPGAs, VPUs (like NCS)

**Using the IE**
- Library: `openvino.inference_engine`
- classes
  - `IE Core`
  - `IE Network`
- methods: `IE Core()`, `IENetwork()`, `iecore.add_extension`, `iecore.query_network(...)`, `net.layers.keys()`, `iecore.load_network`

**IE requests**
- Synchronous → `infer(..)`
- Asynchronous → `start_async(...)`, `wait(..)`

**Handle Results**
- `inputs` → the image frame
- `outputs` → the results
- `latency` → the inference time of current request