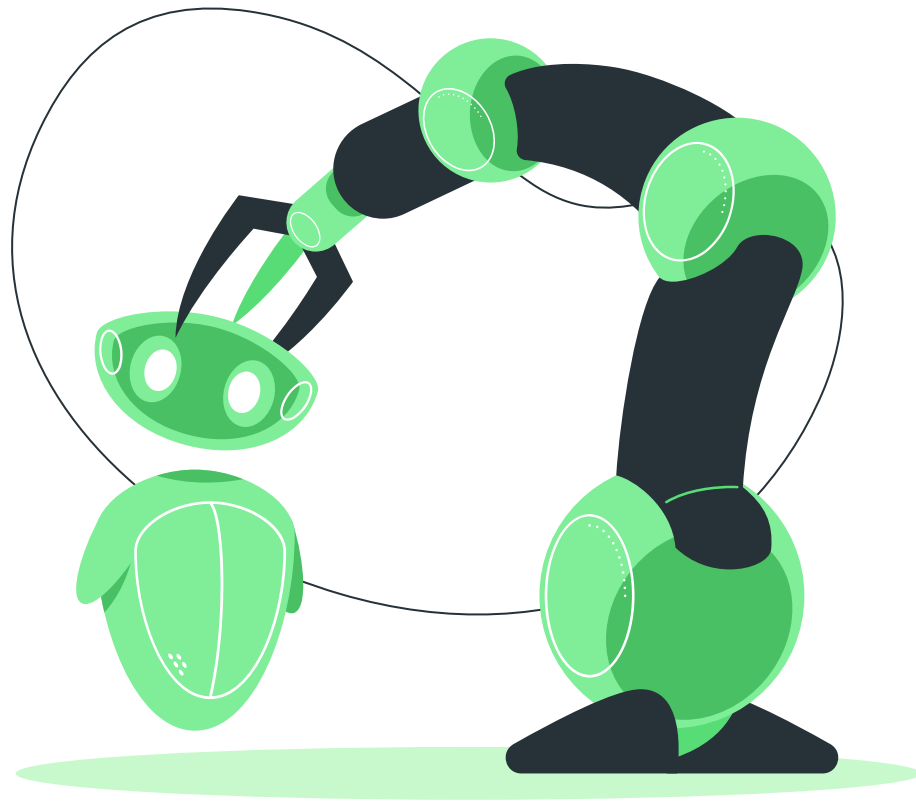


EdgeAI

Evan Galli
Jilian Lubrat
Eliot Menoret
Antoine-Marie Michelozzi



Sommaire

1

Contexte

Objectifs du projet et contraintes de l'IA embarquée.

2

Plateformes

Matériels testés et capacités de calcul.

3

Protocole

Méthodologie de mesure des performances.

4

Optimisations

Techniques d'accélération du modèle IA.

5

Résultats

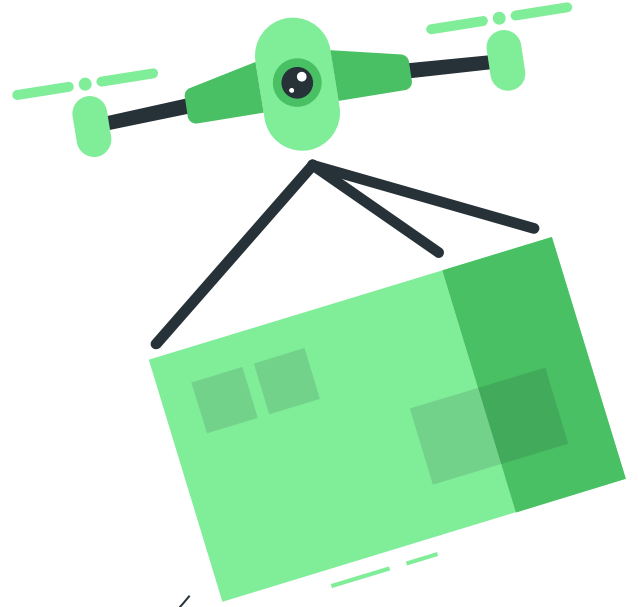
Comparaison des performances obtenues.

6

Conclusion

Bilan et choix optimaux selon l'usage.

1. Contexte



Contexte

Projet de robot suiveur de personne via une caméra, nécessitant un modèle de segmentation

Objectif : Analyser les performances du modèle **YOLOv11n-seg** avec des ressources limitées

Enjeu de l'étude : Comparer les performances sur différents supports matériels de **différentes optimisations** du modèle afin de déterminer le meilleur **compromis** modèle + matériel



Plateformes



AMD Ryzen 5 8645HS

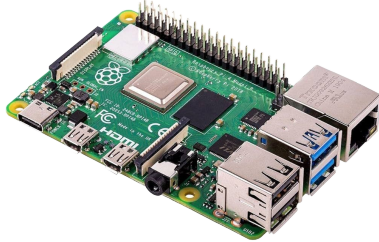
- 6 cœurs / 12 threads
- Fréquence : **4.3 GHz**
- RAM (sur laptop) **32 Go** DDR5



NVIDIA 4060 **Laptop**

- 3072 CUDA cores
- Fréquence GPU : **2.1 GHz**
- Mémoire : **8 Go** GDDR6

Plateformes



Raspberry Pi 4

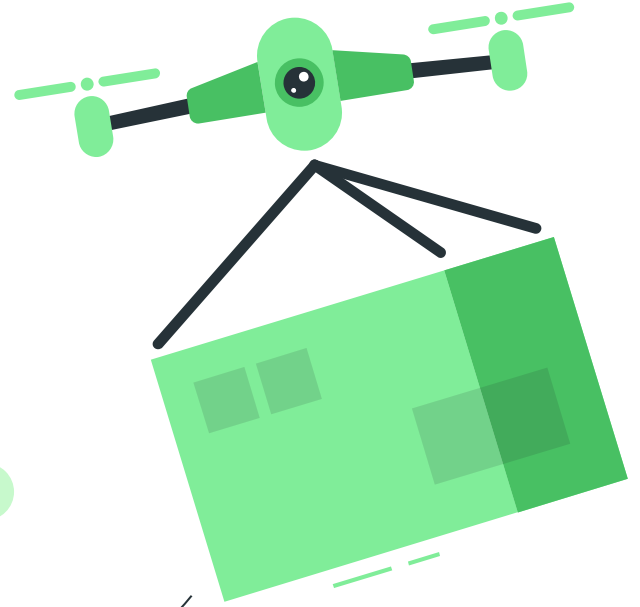
- CPU Broadcom BCM2711
- 4 cœurs
- ARM Cortex-A72 **1.5 GHz**



Caméra Luxonis Oak-D Pro V2

- VPU Intel Movidius Myriad X
- 16 cœurs SHAVES
- Fréquence : **~700 MHz**

3. Protocole de benchmark



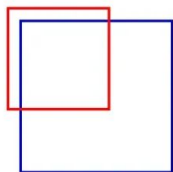
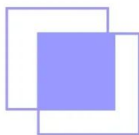
Métriques sélectionnées

Performance Temporelle	Efficacité Énergétique	Qualité de Prédiction
Objectif : Maximiser les images par seconde pour un suivi de personne fluide et sans saccades.	Objectif : Minimiser la consommation énergétique pour utilisation embarquée.	Objectif : Maximiser la précision pour s'assurer que les personnes soient bien détectées

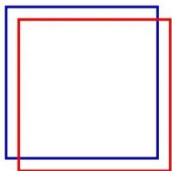
mAP (mean Average Precision)

Pour déterminer si vrai positif :

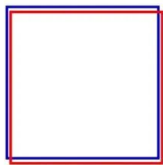
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



Poor



Good



Excellent

$$\text{Précision} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Rappel} = \text{TP} / (\text{TP} + \text{FN})$$

**AP = Aire sous la courbe
Précision-Rappel**

$$mAP = \frac{1}{n} \sum_{k=1}^{k=n} AP_k$$

Calculé pour différent seuil de l'IoU
[0.5;0.95] puis moyenné

Pipeline benchmark

Dataset

- Utilisation du dataset **coco-person** sur **Kaggle**

Objectif

- Comparer les performances uniquement sur la classe **person**



Images extraites du dataset coco-person

Pipeline benchmark

Fonctionnement Pipeline



Itère sur les
images du
dataset



Appelle une
interface
d'inférence unique
(**InferenceBackend**)



Mesure du temps
d'inférence (**avg,
min, max, std**)

Evalue mAP +
mAP mask via
COCOeval



Export des
résultats en
CSV

Inference Backend

Role

- Classe abstraite **standardisant l'inférence**
- Garantit des résultats comparables, indépendamment du **matériel** ou du **framework**

Pré-traitement

- Redimensionnement
- Letterboxing

Post-traitement

- Décodage YOLO
- Application des **masques** et **bounding boxes**
- **Retrait du letterboxing** pour reprojection sur l'image originale

Implémentations

OnnxBackend

- Sert à exécuter l'inférence YOLO sur **PC/RPi** via **ONNX Runtime**
- Permet d'évaluer **CPU** et **GPU** avec la même interface

OakDBackend

- Sert à exécuter l'inférence sur l'**OAK-D Pro** (NPU intégré)
- **Envoie l'image** à la caméra avec la lib **Depthai**, récupère les sorties du modèle

Consommation

Méthodologie de mesure

- Mesure basée sur le **delta de consommation** entre l'état idle et l'inférence
- Permet une comparaison **équitable** avec la caméra embarquée, qui n'exécute pas de système d'exploitation

Environnement contrôlé

- **Nettoyage** du système hôte
- Arrêt des processus et services **non essentiels**
- Réduction maximale des tâches en arrière-plan pour limiter le **bruit de mesure**

Consommation

GPU (NVIDIA)

- Utilisation de **pynvml** wrapper de **NVML**, lib officielle NVIDIA pour monitorer et gérer les GPU

CPU

- Enregistrement des consommations via le logiciel **HWiNFO** (via export CSV)

Raspberry/Oak-d Pro

- Enregistrement vidéo du dongle USB mesurant la consommation
- Récupération d'une image par seconde
- Envois à **Gemini** des images pour récupérer les valeurs



Quantization

Principe

- Réduction de la **précision** des poids du modèle

Objectifs

- **Accélérer** l'inférence
- Réduire la consommation
- **Diminution** l'empreinte mémoire

Formats utilisés

- **FP16** : supporté sur toutes les plateformes
- **INT8** : Raspberry Pi, CPU & GPU (non supporté sur **Oak-D Pro V2**)
- En dessous perte de précision trop importante

Quantization

Statique

- Quantization effectuée **avant l'inférence**
- **Nécessite** un jeu de données de **calibration**
- **Poids** et **activations** quantifiés
- Temps d'inférence / Taille du modèle **réduits**

Dynamique

- **Poids** quantifiés hors ligne
- **Activations** quantifiées à l'exécution
- **Réduction** taille moins importantes que **statique**
- Gains matériels **limités** (scale pour les activations)
- **Plus simple:** Pas de ré-entraînement

Pruning

Structuré

- **Suppression définitive** de poids ou de canaux
- Modèle **plus léger**
- **Accélération réelle** à l'inférence
- **Ré-entraînement** nécessaire

Non structuré

- **Poids ignorés** à l'exécution
- **Taille** du modèle **inchangée**
- **Aucun gain** matériel réel
- Pas de ré-entraînement

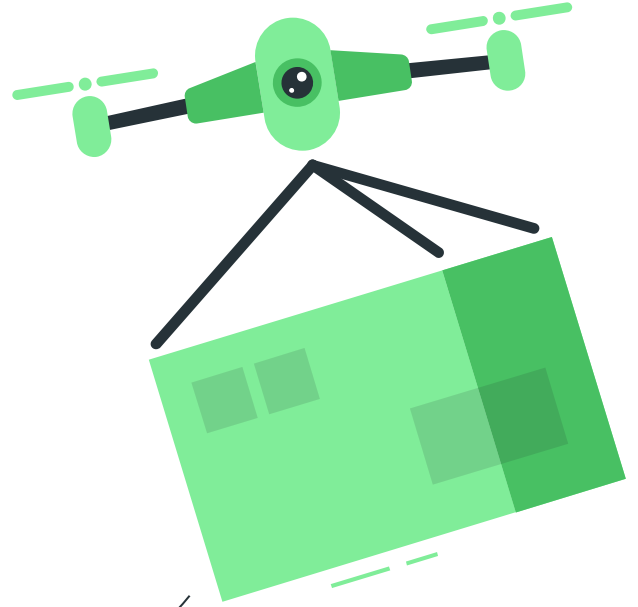
Pruning structuré

- Utilisation du projet YOLO-Pruning-RKNN par heyongxin233
- Pruning structurel qui supprime des canaux ou des couches entières. Permet d'accélérer réellement l'inférence.
- Adaptation du projet pour utiliser le dataset **coco-person** durant le réentraînement

Optimisations par plateformes

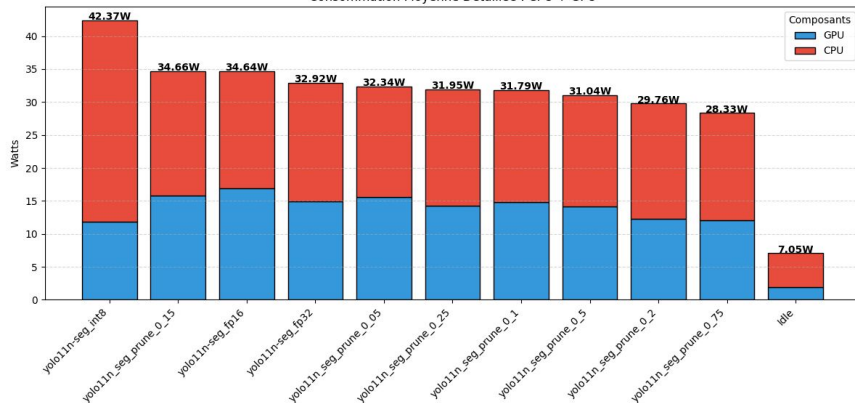
OAK-D Pro	Raspberry Pi 4	PC
<ul style="list-style-type: none">• Allocation de 8 SHAVES (conseillé par la documentation officielle)	<ul style="list-style-type: none">• Paramétrage de ONNX Runtime :<ul style="list-style-type: none">○ Graph Optimization○ Multi-threading○ Gestion RAM○ Mode Séquentiel	<ul style="list-style-type: none">• Utilisation du GPU pour paralléliser les calculs dans le cas du benchmark CPU + GPU

5. Résultats

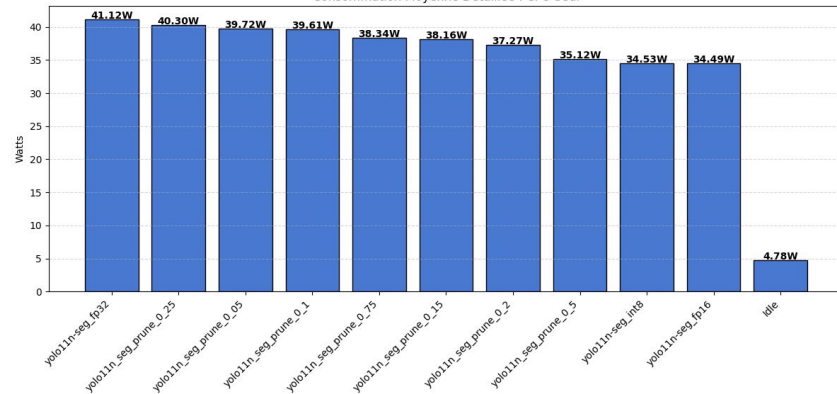


Résultats

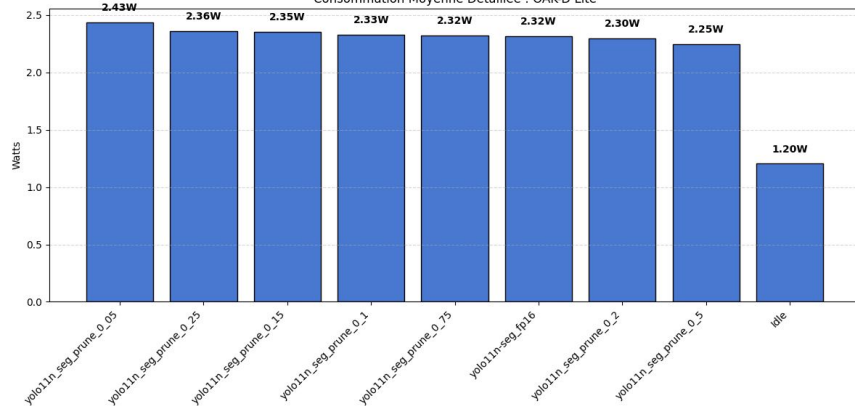
Consommation Moyenne Détaillée : CPU + GPU



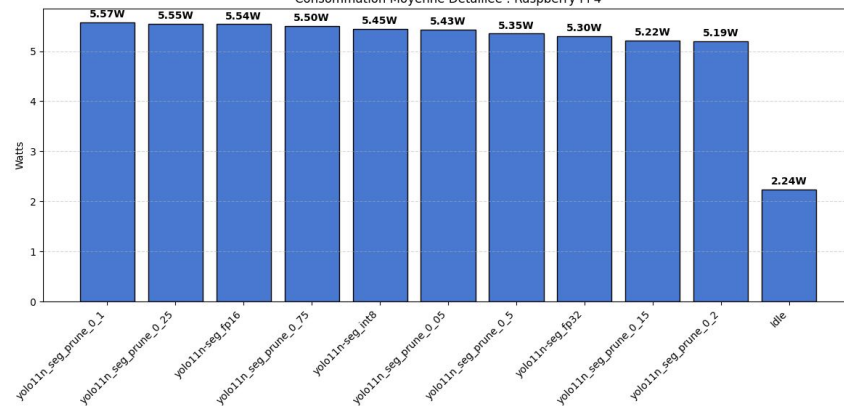
Consommation Moyenne Détaillée : CPU Seul



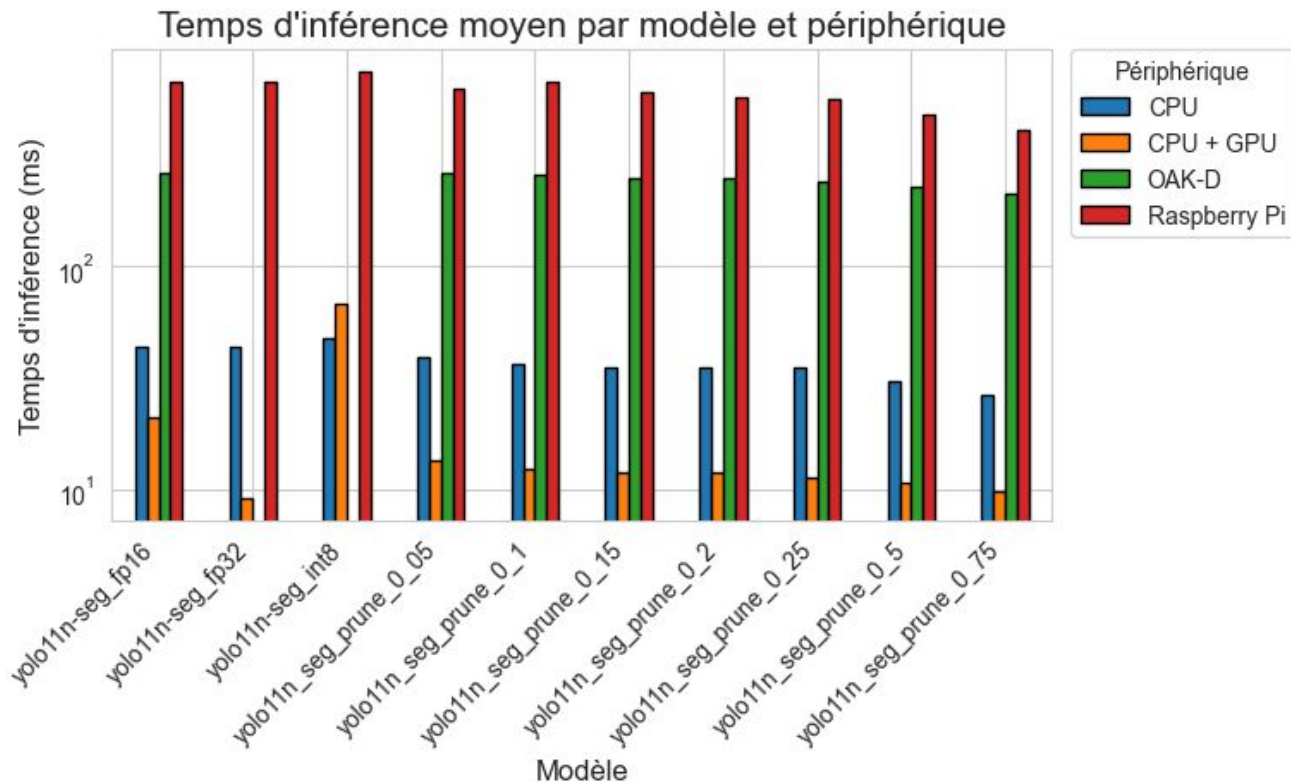
Consommation Moyenne Détaillée : OAK-D Lite



Consommation Moyenne Détaillée : Raspberry Pi 4

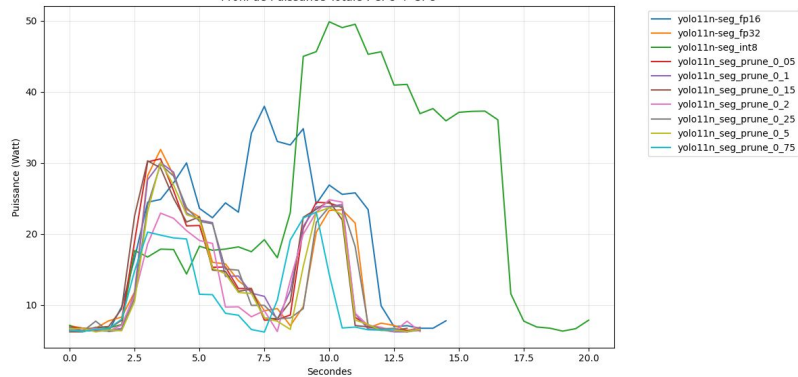


Résultats

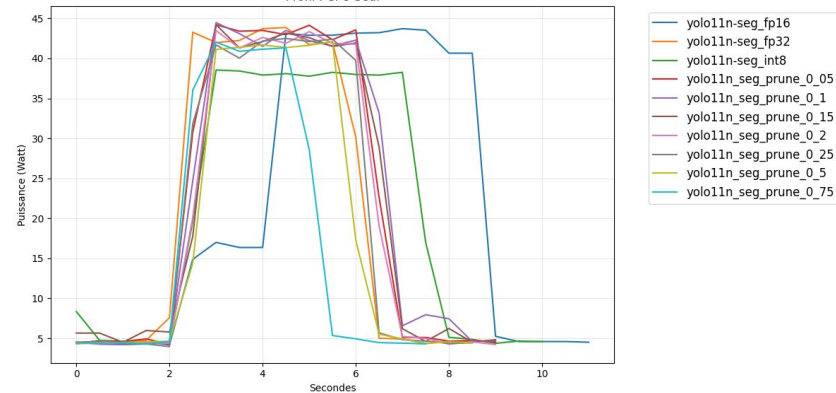


Résultats

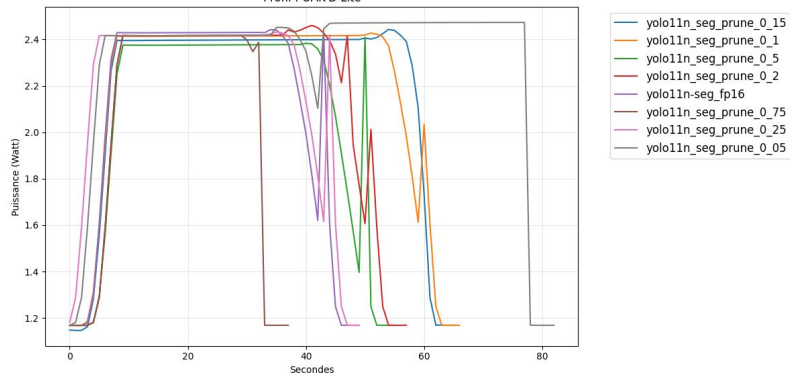
Profil de Puissance Totale : CPU + GPU



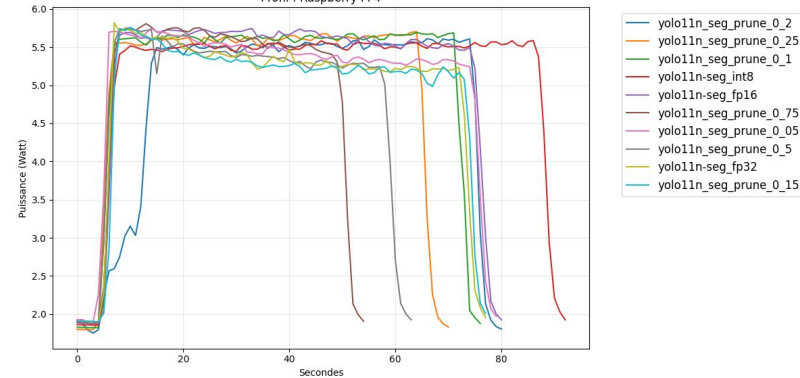
Profil : CPU Seul



Profil : OAK-D Lite

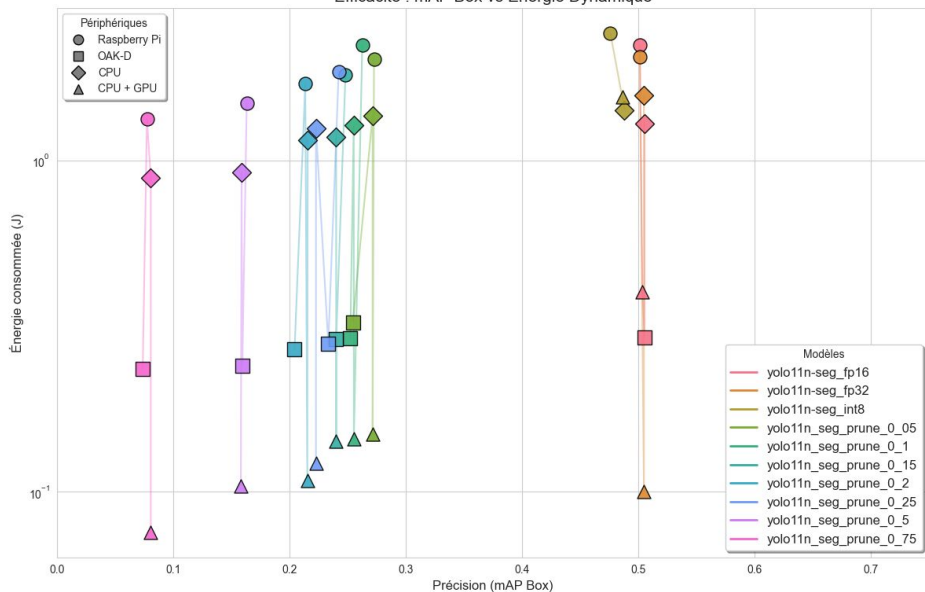


Profil : Raspberry Pi 4

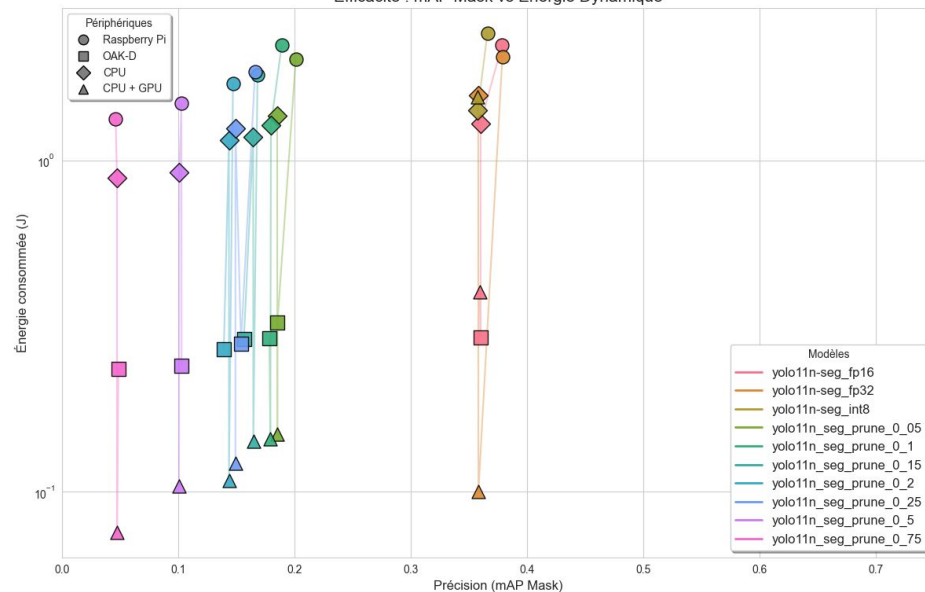


Résultats

Efficacité : mAP Box vs Énergie Dynamique



Efficacité : mAP Mask vs Énergie Dynamique





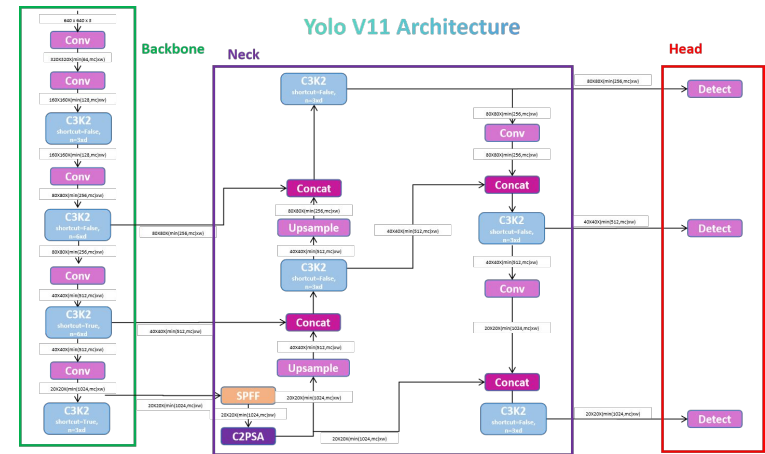
Choix Optimal : Device

- La caméra **OAK-D Pro** est le choix optimal
- Sa consommation énergétique est **légèrement inférieure** à CPU+GPU
- Le mAP (Mean Average Precision) obtenu est **similaire** sur les autres plateformes
- L'utilisation sur batterie permet une solution **embarquée**
- Meilleur compromis entre modèle et matériel



Choix Optimal : Modèle

- Le modèle en précision **FP16** est le meilleur choix à utiliser avec la caméra
- La consommation d'énergie est similaire aux autres modèles
- **mAP bien meilleur** que pour les autres modèles prunés





Merci

Avez-vous des questions ?