

XiNet

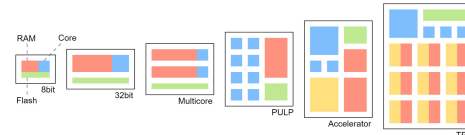
Challenging the **efficiency** of depthwise convolutions for **edge and tinyML**:

- Novel convolutional block optimizing **latency** and **energy** usage
- **Benchmarking** on multiple embedded platforms
- **Hardware Aware Scaling**: from hardware constraints to neural architecture

Hardware-Aware Scaling

Three main computational constraints in different embedded devices:

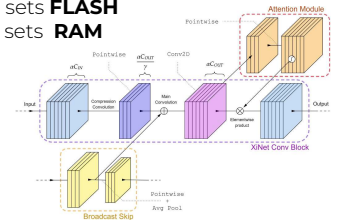
- **FLASH**: stores network parameters
- **RAM**: stores intermediate tensors
- **MAC/s**: determines latency & energy



Convolutional Block

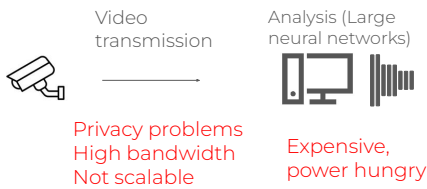
Designed from **real world** efficiency measurements on various platforms. Three hyperparameters:

- α : sets **MAC**
- β : sets **FLASH**
- γ : sets **RAM**



Object/Pose Detection

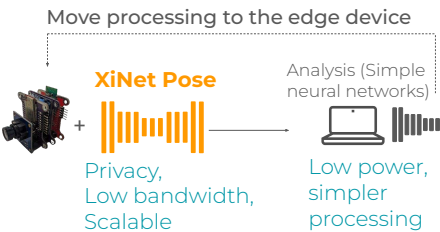
Typical pipeline:



Reducing bandwidth by 3 orders of magnitude!



On-device:



Results

Networks scaled using **Hardware Aware Scaling**:

- MCU**: STM32 - 100MMAC/s, 2MB Flash, 1MB Ram
- TPU**: K210 - 1GMAC/s, 16MB Flash, 5MB Ram
- MPU**: rPi 4B - 16GMAC/s, 16GB SD, 4GB Ram



Raspberry Pi 4B:

- $\alpha=1.0$ $\beta=1.0$ $\gamma=4.0$
- Speed: **63.6 fps**
- Power: **14.89 W**



Kendryte K210:

- $\alpha=0.75$ $\beta=1.0$ $\gamma=4.0$
- Speed: **31.7 fps**
- Power: **410 mW**

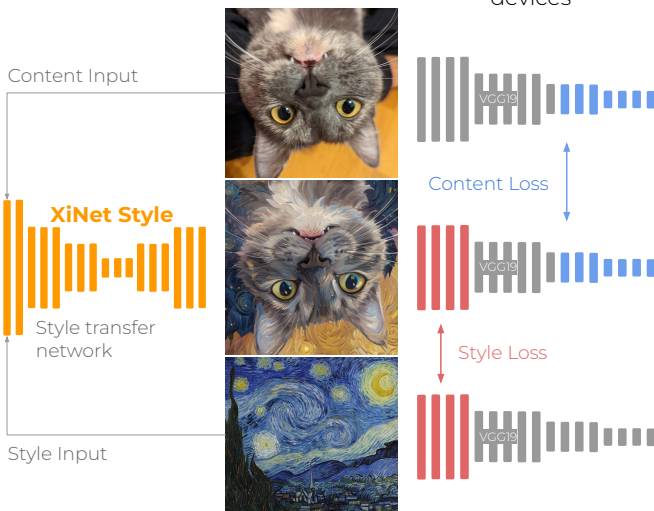


ST STM32H743:

- $\alpha=0.33$ $\beta=0.8$ $\gamma=5.0$
- Speed: **5.5 fps**
- Energy: **72.4 mW**

Style transfer

Efficient **image generation** on edge devices

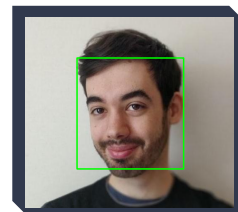


One shot image generation can be used for **anonymization** while preserving semantic content - removing personal information for **downstream tasks**

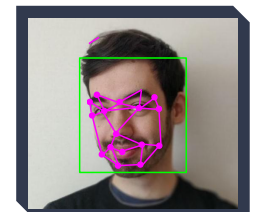
XiSwap

Single target face swapping

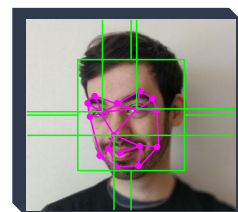
in 4 steps using 3 networks:



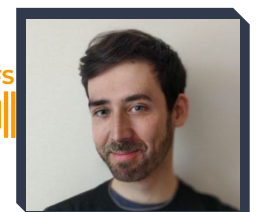
1. Face Detection
XiNet + Yolo



2. Landmark Detection
XiNet + PFLD



3. Face Alignment
XiNet + PFLD



4. Face Generation
XiNet GAN

