SkyNet: a Hardware-Efficient Method for Object Detection and Tracking on Embedded Systems

Xiaofan Zhang¹, Haoming Lu¹, Cong Hao¹, Jiachen Li¹, Bowen Cheng¹, Yuhong Li¹, Kyle Rupnow², Jinjun Xiong^{3,1}, Thomas Huang¹, Honghui Shi^{3,1}, Wen-mei Hwu¹, Deming Chen^{1,2}

¹C³SR, UIUC ²Inspirit IoT, Inc ³IBM Research





Outline:

1) Background & Challenges

Edge AI is necessary but challenging

2) Motivations

Two major issues prevent better AI quality on embedded systems

3) The Proposed SkyNet Solution

A bottom-up approach for building hardware-efficient DNNs

4) Demonstrations on Object Detection and Tracking Tasks

Double champions in an international system design competition Faster and better results than trackers with ResNet-50 backbone 5) Conclusions

Cloud solutions for AI deployment

Major requirements:

- High throughput performance
- Short tail latency





Language Translation



Voice-activated assistant

Recommendations

Google 发布了一款屏保 App,自动更新最精美的

法冬新款里



Video analysis

Why still need Edge solutions?

Communication

Privacy



Demanding AI applications cause great challenges for Edge solutions.

We summarize three major challenges

Latency

Edge AI Challenge #1 Huge compute demands



https://openai.com/blog/ai-and-compute/

Edge AI Challenge #1 Huge compute demands



https://openai.com/blog/ai-and-compute/

Edge AI Challenge #2 Massive memory footprint



18]

7

Edge AI Challenge #2 Massive memory footprint

- HD inputs for real-life applications
 - 1) Larger memory space required for input feature maps
 - 2) Longer inference latency
- Harder for edge-devices
 - 1) Small on-chip memory
 - 2) Limited external memory access bandwidth



Edge AI Challenge #3 Real-time requirement

- Video/audio streaming I/O
 - 1) Need to deliver high throughput
 - 24FPS, 30FPS ...



Edge AI Challenge #3 Real-time requirement

- Video/audio streaming I/O
 - 1) Need to deliver high throughput
 - 24FPS, 30FPS ...
 - 2) Need to work for real-time
 - E.g., millisecond-scale response for self-driving cars, UAVs
 - Can't wait for assembling frames into a batch

Outline:

1) Background & Challenges

Edge AI is necessary but challenging

2) Motivations

Two major issues prevent better AI quality on embedded systems

3) The Proposed SkyNet Solution

A bottom-up approach for building hardware-efficient DNNs

4) Demonstrations on Object Detection and Tracking Tasks

Double champions in an international system design competition

Faster and better results than trackers with ResNet-50 backbone

5) Conclusions

A Common flow to design DNNs for embedded systems

Various key metrics: Accuracy; Latency; Throughput; Energy/Power; Hardware cost, etc.

Trained DNNs	SW-related	HW-related	Implementation on
	Optimization	Optimization	embedded devices
 More focus on accuracy Excessively complicated for IoT <i>Step 1</i> 	 Quantization Pruning Layer fusion Conv variation <i>Step 2</i> 	 Parallel factors adjustment Resource allocation I/O optimizations <i>Step 3</i> 	Step 4

It is a top-down flow: form reference DNNs to optimized DNNs

Object detection design for embedded GPUs

Target NVIDIA TX2 GPU

~665 GFLOPS @1300MHz

Input resizing <a>2 Pruning <a>3 Quantization <a>4 TensorRT
 Multithreading

GPU-Track	Reference	Software Optimizations	Hardware Optimizations
'19 2 nd Thinker	ShuffleNet + RetinaNet	123	(5)
'19 3 rd DeepZS	Tiny YOLO	-	5
'18 1 st ICT-CAS	Tiny YOLO	1234	-
'18 2 nd DeepZ	Tiny YOLO	-	5
'18 3 rd SDU-Legend	YOLOv2	123	(5)

[From the winning entries of DAC-SDC'18 and '19]

Object detection design for embedded FPGAs

Target Ultra96 FPGA

~144 GFLOPS @200MHz



(1) Input resizing (2) Pruning (3) Quantization

5 CPU-FPGA task partition 6 double-pumped DSP 7 pipeline 8 clock gating

FPGA-Track	Reference	Software Optimizations	Hardware Optimizations
'19 2 nd XJTU Tripler	ShuffleNetV2	23	568
'19 3 rd SystemsETHZ	SqueezeNet	123	7
'18 1 st TGIIF	SSD	123	56
'18 2 nd SystemsETHZ	SqueezeNet	123	7
'18 3 rd iSmart2	MobileNet	123	57

[From the winning entries of DAC-SDC'18 and '19]

Drawbacks of the top-down flow

1) Hard to balance the sensitivities of DNN designs on software and hardware metrics



- 2) Difficult to select appropriate reference DNNs at the beginning
 - Choose by experience
 - Performance on published datasets

Outline:

1) Background & Challenges

Edge AI is necessary but challenging

2) Motivations

Two major issues prevent better AI quality on embedded systems

3) The Proposed SkyNet Solution

A bottom-up approach for building hardware-efficient DNNs

4) Demonstrations on Object Detection and Tracking Tasks

Double champions in an international system design competition

Faster and better results than trackers with ResNet-50 backbone

5) Conclusions

The proposed flow

To overcome drawbacks, we propose a bottom-up DNN design flow:

- No reference DNNs; Start from scratch;
- Consider HW constraints; Reflect SW variations



The proposed flow [overview]

It is a three-stage flow

Select Bundles -> Explore network architectures -> Add features



The proposed flow [stage 1]

Start building DNNs from choosing the HW-aware Bundles
Goal: Let Bundles capture HW features and accuracy potentials



- Prepare DNN components
- Enumerate Bundles
- Evaluate Bundles (Latency-Accuracy)
- Select those in the Pareto curve

The proposed flow [stage 2]

Start exploring DNN architecture to meet HW-SW metrics
 Goal: Solve the multi-objective optimization problem



- Stack the selected Bundle
- Explore two hyper parameters using *PSO* (channel expansion factor & pooling spot)
- Evaluate DNN candidates (Latency-Accuracy)
- Select candidates in the Pareto curve

- Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



Fitness Score:

$$Fit_j^i = Acc_j^i + \alpha \cdot (Est(n_j^i) - Tar)$$

Candidate accuracy

Candidate latency in hardware

Targeted latency

 $lpha\,$ factor to balance accuracy and latency

- > Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



- > Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



- Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



- Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



- > Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



- Adopt a group-based **PSO** (particle swarm optimization)
- Multi-objective optimization: *Latency-Accuracy*
- Group-based evolve: Candidates with the same **Bundle** are in the same group



The proposed flow [stage 3]

Add more features if HW constraints allow

Goal: better fit in the customized scenario



- For small object detection, we add feature map bypass
- Feature map reordering
- For better HW efficiency, we use ReLU6

The proposed flow [HW deployment]

> We start from a well-defined accelerator architecture

Two-level memory hierarchy to fully utilize given memory resources IP-based scalable process engines to fully utilize computation resources



The proposed flow [HW deployment]

> We start from a well-defined accelerator architecture



Enable fast performance evaluation



Outline:

1) Background & Challenges

Edge AI is necessary but challenging

2) Motivations

Two major issues prevent better AI quality on embedded systems

3) The Proposed SkyNet Solution

A bottom-up approach for building hardware-efficient DNNs

4) Demonstrations on Object Detection and Tracking Tasks

Double champions in an international system design competition

Faster and better results than trackers with ResNet-50 backbone

5) Conclusions

Demo #1: an object detection task for drones

System Design Contest for *low power object detection* in the IEEE/ACM Design Automation Conference (DAC-SDC)



- DAC-SDC targets single object detection for real-life UAV applications Images contain 95 categories of targeted objects (most of them are small)
- > Comprehensive evaluation: *accuracy, throughput, and energy consumption*

 $TS_i = R_{IoU_i} \times (1 + ES_i)$

Demo #1: DAC-SDC dataset

The distribution of target relative size compared to input image

31% targets < 1% of the input size 91% targets < 9% of the input size





Demo #1: the proposed DNN architecture



- 13 CONV with 0.4 million parameters
- For Embedded FPGA: Quantization, Batch, Tiling, Task partitioning
- For Embedded GPU: Task partitioning



Demo #1: Results from DAC-SDC [GPU]

Evaluated by 50k images in the official test set





Designs using TX2 GPU

Demo #1: Results from DAC-SDC [FPGA]

Evaluated by 50k images in the official test set



'19 Designs using Ultra96 FPGA



'18 Designs using Pynq-Z1 FPGA



Demo #2: generic object tracking in the wild

- We extend SkyNet to real-time tracking problems
- > We use a large-scale high-diversity benchmark called *Got-10K*
 - *Large-scale:* 10K video segments with 1.5 million labeled bounding boxes
 - Generic: 560+ classes and 80+ motion patterns (better coverage than others)



[From Got-10K]

Demo #2: Results from Got-10K

Evaluated using two state-of-the-art trackers with single 1080Ti

Backbone	AO	$SR_{0.50}$	$SR_{0.75}$	FPS	
AlexNet	0.354	0.385	0.101	52.36	
ResNet-50	0.365	0.411	0.115	25.90	Siı
SkyNet	0.364	0.391	0.116	41.22	VS

SiamRPN++ with different backbones

Similar AO, 1.6X faster vs. ResNet-50

SiamMask with different backbones

Backbone	$\mid AO$	$SR_{0.50}$	$ SR_{0.75} $	$\mid FPS$	
ResNet-50	0.380	0.439	0.153	17.44	Slightly better AO, 1.7X faster
SkyNet		0.442	0.158	30.15	vs. ResNet-50

Outline:

1) Background & Challenges

Edge AI is necessary but challenging

2) Motivations

Two major issues prevent better AI quality on embedded systems

3) The Proposed SkyNet Solution

A bottom-up approach for building hardware-efficient DNNs

4) Demonstrations on Object Detection and Tracking Tasks

Double champions in an international system design competition

Faster and better results than trackers with ResNet-50 backbone

5) Conclusions

Conclusions

- We presented SkyNet & a hardware-efficient DNN design method
- a bottom-up DNN design flow for embedded systems
- an effective way to capture realistic HW constraints
- a solution to satisfy demanding HW and SW metrics
- SkyNet has been demonstrated by object detection and tracking tasks
- Won the double champions in DAC-SDC
- Achieved faster and better results than trackers using ResNet-50



- Scan for paper, slides, poster, code, & demo
- Please come to Poster #11

Conference on Machine Learning and Systems (MLSys) 2020



Thank you

