

ECOLE NATIONALE SUPÉRIEURE D'INFORMATIQUE



HyT-NAS: Hybrid Transformers Neural Architecture Search for Edge Devices

Date: 13/10/2022

MECHARBAT Lotfi Abdelkrim, BENMEZIANE Hadjer, OUARNOUGHI Hamza, NIAR Smail







Confidence

Bounding Box

High accuracy in various fields, including object recognition.

Extremely flexible due to the wide variety of hyperparameters that control them.

Edge Al

- Unreliable (depends on network quality).
- Slow process for real-time applications.
- Not suitable for critical applications.

Gigantic architectures, models are too big to fit in Edge devices.

Huge computational complexity, not fast enough for inference in Edge

High power consumption, drains the limited power source (battery) of Edge devices.

More reliable.
No data transfer over the network.
Preserve confidentiality.

Hyperparameters optimization

Why is it difficult to make a good choice of hyperparameters manually?

The space of possible configurations is of immense size

lf

High cost evaluation which consists in training deep learning architectures

Ex: The learning time of ViT on ImageNet1k for 100 epochs on 8 NVIDIA A100- 40GB GPUs is 65 hours. source: https://ai.facebook.com/blog/significantly-faster-vision-transformer-training/

Size of the space of possible configurations= 1.099 * 10^12 With 1s/eval, the exploration of this space requires more than 30,000 years Objective

Propose an efficient hardware-aware neural architecture search method to find Hybrid **Transformer models** that are fast, deployable on small edge devices and effective for Visual **Object Recognition.**

Study Case

Image Classification

Object Detection

State of the art

Hardware Aware Neural Architecture Search (HW-NAS)

State of the art

Hardware Aware Neural Architecture Search (HW-NAS)

Hybrid Search Space

Propose an Initial search space : Accuracy-focused study SOTA architectures for Visual Object Recognition.

- Too big to efficiently explore $\sim 10^{27}$
- Does not consider hardware constraints

Efficiency analysis : Comparative study of the efficiency of SOTA models and operations on edge devices according to hardware metrics such as Latency, Memory consumption, Size and Throughput.

Hybrid models are more likely to be deployed on edge devices.

Hyperparameters such as the number of heads and the embedding size have more impact on the size and efficiency of attention blocks than others.

Seconde

Hybride Search Space

Description

Latency (s)

Block	Hyperparameter	Values		
	Number of blocks	[1, 2, 3, 4]		
onvolution Block	Expand ratio	[1x, 2x, 4x]		
	Out channel size	[8, 16, 24, 32]		
	Expand ratio	[1x, 2x, 4x]		
	Channel size	[1x, 1.5x, 2x]		
Attention Block	Number of heads	[1, 2, 4]		
	Feed forward ratio	[1x, 1.5x, 2x]		

Accuracy (%)

Search Strategy

$max_{\alpha \in HySS}$ Accuracy(α), Throughput(α) subject to N parameters(α) $\leq MaxN parameters$

all objectives.

Search Strategy Study

Surrogate	 XgBoost, XgBRanker Feed Forward Networks (FFN) Gaussian Process (GP) Bayesian Neural Network (BNN) 	Method	Surrogate	Acquisition function	Multi-objective solver	Selection method	Performance (Avg Number of discovered paretos)
		Random					3.68/14
		CMA-ES					5.45/14
Acquisition	 UCB (Upper Confidence Bound) EI (Expected Improvement) 	NSGAII					6.06/14
		MOBO std	GP	EI	NSGAII	None	5.4/14
Multi-objective solver	• NSGAII	HyT-Search	BNN	EI	NSGAII	HVI	5.2/14
		HyT-Search	FFN (1layer)	EI	NSGAII	HVI	10.2/14
		HyT-Search	FFN(2layer)	UCB	NSGAII	Random	11.4/14
Selection method	• HVI (Hypervolume Improvement)	HyT-Search	XGBoost	UCB	NSGAII	Dominance	12.6/14
	RandomDominance	HyT-Search	XgBoost	UCB	NSGAII	HVI	13.7/14

Benchmark: Reproducible and Efficient Benchmarks for Hyperparameter Optimization (<u>https://github.com/Este1le/hpo_nmt</u>)

Search Strategy Evaluation

Evaluation strategy

Results

Visual Wake Words

HyT-NAS-BL outperforms MobileVit variants while significantly reducing latency and the number of parameters.

HyT-NAS-BA is largely more accurate with lower latency than all the others and a smaller size than the most.

HyT-NAS-O outperforms the 90% in accuracy with a latency and size more optimal than all the others.

Person Detection

	MobileNetV3
Our HyT-NAS_BO detector achieves better accuracy than mobilenetV3 while being much smaller (more than 5x).	MobileViT-XS
Our HyT-NAS_BO detector achieves similar accuracy as MobileViT-XXS while being	MobileViT-XXS
	HyT-NAS_BO

50

• Accuracy (mAP) • Size (100k)

Propose a new method of automatic search of neural architecture adapted to the hardware called "HyT-NAS".

Realize a comprehensive study of Vision Transformers models for visual object recognition on several hardware platforms.

Propose a new hybrid search space that includes convolution and attention blocks targeting small edge devices.

Propose a new search strategy aims to accelerate convergence by finding good architectures in a relatively small number of evaluations.

Perspective

14

Expanded the search space by allowing interchanging of attention and convolution blocks

Consider other metrics in the optimization such as: energy consumption.

Add semantic segmentation as a use case.

S

Thank you for your attention

HyT-NAS: Hybrid Transformers Neural Architecture Search for Edge Devices

- MECHARBAT Lotfi Abdelkrim (ESI ex INI)
 - hl_mecharbat@esi.dz
- https://www.linkedin.com/in/lotfi-abdelkrim-mecharbat-7740b3164/
 - https://github.com/meclotfi/HyT-NAS-Search-Algorithm