

Hardware-aware Neural Architecture Search

Annika Österdiekhoff

Embedded Systems

University Duisburg-Essen

Duisburg, Germany

annika.oesterdiekhoff@stud.uni-due.de

Abstract—Neural Architecture Search is a grown topic for automatically design deep learning models. With the raising of IoT and mobile devices nowadays, deep learning models have to be fit on specific hardware. Therefore, hardware-aware NAS is developed which does not only consider to maximize the accuracy of a deep learning model but also the hardware usability measured by FLOPs, model size, latency, energy consumption, etc. This paper gives a brief overview about hardware-aware neural architecture search by focusing on the problem definition of hardware-aware NAS, its search space and search strategy and its challenges and limitations.

Index Terms—Hardware-aware Neural Architecture Search, Optimization, Hardware Platforms

I. INTRODUCTION

Deep learning models are gaining more and more interest in various subjects for example in image recognition [2] or Natural Language Processing [3]. These models consist of different layers and parameters. Layers are operators like convolution or activation whereas parameters also known as hyper-parameters are pre-defined properties of the architecture or the training algorithm. Architecture Parameters are for example the stride and filter of an convolution layer whereas a training parameter can be the number or epochs.

There exist a huge number of building up this deep learning model out of the layers and parameters. So, it is difficult for the developer to design the number and type of nodes and the connection between them because of the multiple differences in data types, tasks and hardware platforms. Therefore, most of the deep learning models are created by hand with multiple experiments or are a variation of already known models which work well. This process of creating and designing a model is very time-consuming and costly, therefore techniques are created in the last years to automate the designing process of a deep learning model. It is called Neural Architecture Search (NAS). The focus of this paper is a subgroup of NAS, the hardware or platform-aware NAS (HW-NAS). The goal of HW-NAS is to use NAS for designing a deep learning model but with respect to optimize the deep learning model for a hardware device. So, for giving an overview about hardware-aware NAS we summarize a recent survey from 2021 [1].

In the following we define a general NAS process in Section II. In Section III possible ways to optimize a deep learning model are listed. We distinguish between single-objective and multi-objective optimization in Section IV. After defining the goals of hardware-aware NAS in Section V, possible hardware platforms are listed in Section VI. The architecture

and hardware search space are explained in Section VII. We go into detail with the accuracy evaluation method, the hardware cost evaluation method and the search algorithms in Section VIII. The NAS frameworks and especially the Microsoft NNI Framework are shortly presented in Section IX. In the end we define the challenges and limitations in Section X by speaking about benchmarking, transferability of AI models and transferability across different hardware platforms.

II. NEURAL ARCHITECTURE SEARCH DEFINITION

A neural architecture process firstly includes a **search space** which consists of possible operators and its connections which create an architecture. The **search strategy** explores this search space and searches for possible architecture candidates. Then the **evaluation methodology** is used to evaluate the accuracy of each model. The evaluation methodology trains the architecture. The models with a high accuracy help to redefine the search space. The whole NAS process is shown in Figure 1 [1]. This training and in general the number of architecture in forms of layers and parameters cause time consumption and large memory footprints. This leads to a lack of using NAS for constrained hardware or in real-time. Therefore, more and more research is focused on solving this issue by using NAS and developing hardware-aware NAS.

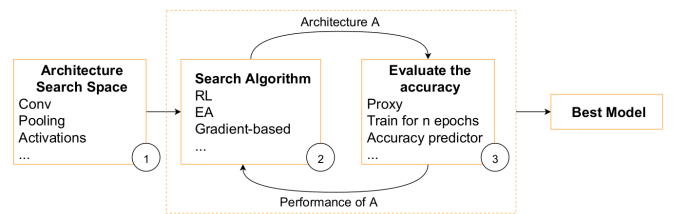


Fig. 1. NAS Process [1]

III. MODEL OPTIMIZATIONS

For solving the drawbacks mentioned above of getting large deep learning models and long training processes, there are some techniques to optimize it.

Firstly you can **compress a model** by decreasing the model size but gaining the same accuracy. Compressing a model can be done for example by *compact the model* which changes the standard operations to more flexible and simpler operations. In addition, one can *decompose the tensor* which means we

shrink the tensors which reduces the size of the deep learning model. What we will mainly focus on for compressing a model is *quantization*. Quantization means that we convert floating point weights and activation's into smaller integers, ideally binaries. This change of representation can be different in each layer, this is called mixed-precision quantization. Another option for compressing a model is *pruning* where the least important weights or operations are pruned to reduce the model size. The importance of a weight or operation can be the weight itself or has to be learned. There is put lots of effort the last years to automate mixed-precision quantization as well as the pruning.

In addition one can apply **hardware-aware NAS** which uses model compression's and takes the hardware usability into account by searching the model from a set of architectures. A relatively new technique for optimizing a deep learning model is **code transformation** which optimizes the operators of an model for a specific hardware.

Figure 2 [1] shows the input and output of compression a model or applying hardware-aware NAS for optimizing a model.

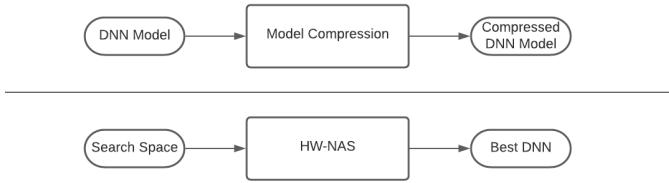


Fig. 2. The Process of optimizing a model by model compression or HW-NAS [1]

IV. HARDWARE-AWARE NAS PROBLEM FORMULATION

In the following we discuss the problem formulation of hardware-aware neural architecture search. In general, neural architecture search is an optimization problem which search for an architecture which maximizes a performance metrics. The performance metric is mostly the accuracy. The formulation is shown in “(1)”.

$$\max_{\alpha \in A} f(\alpha, \delta) \quad (1)$$

So, A is the search space which includes the set of architectures and α is one architecture from this set. δ is the used data set and f is the performance measure, i.e. the accuracy. The formulation says that we try to find the architecture α of the search space A which maximizes the accuracy f for the used data set δ .

By doing hardware-aware neural architecture search one has not only the accuracy metric which should be maximized but also hardware objectives like for example latency. There exist two different problem formulation for hardware-aware neural architecture search, the single and multi-objective optimization.

A. Single-Objective Optimization

Single-objective optimization means that only one objective is maximized like explained above for example the accuracy. For solving the problem of also having hardware objectives one tries to transform the multi-objective optimization into a single-objective optimization. For doing this there exist two methods. **Two-Stage optimization** means that we use the problem formulation of original neural architecture search which maximizes the accuracy and optimize in the second stage the architecture for a specific target. This method often does not work well because the best founded architecture in the first stage is maybe not the best for the specific target hardware. The alternative to this method is **constrained optimization**. Here we again try to find the architecture with the highest accuracy but we define some thresholds for hardware constraints like latency or energy consumption which are not allowed to be exceeded by the founded architecture.

B. Multi-Objective Optimization

The difference to the single-objective optimization which is used in original NAS is multi-objective optimization. The formulation is shown here: “(2)”.

$$\max_{\alpha \in A} f_1(\alpha, \delta), f_2(\alpha, \delta), \dots, f_n(\alpha, \delta) \quad (2)$$

So we now do not only have the accuracy which we have to optimize but various objectives, e.g. latency and energy consumption. These objectives can be conflicting each other for example try to have a minimal search space and gaining the best accuracy. When this occurs we can only try to find the Pareto-optimal solution. The Pareto-optimal solution is a set of architectures where the changing of one measurement, for example improving the accuracy, always leads to a decrease for a different measurement, for example higher latency. **HOW TO CITE?** It is not possible to improve one objective without drawbacks to the other objectives.

There exist two approaches for this. Firstly the **scalarization method** which includes the performance measurements f 's into one aggregation function like a weighted sum or weighted product. This transforms the multi-objective optimization into a single-objective optimization. Unfortunately by having fixed weights in the aggregation function, it could be that not all Pareto solutions are found. For solving this issue one can run multiple runs which on the other hand then again leads to more needed computing power etc. Another method to find the Pareto-optimal solution is defining **heuristics**. By doing this one does not form the multi-objective optimization problem into a single-objective one and therefore get a set of architectures as solutions along the Pareto front.

The Figure 3 is leaned on a figure of the survey [1] we summarize and shows the different possibilities of formulating the NAS problem.

V. GOALS OF HARDWARE-AWARE NAS

There are three categories of goals of hardware-aware neural architecture search. Firstly, the category **single target, fixed**

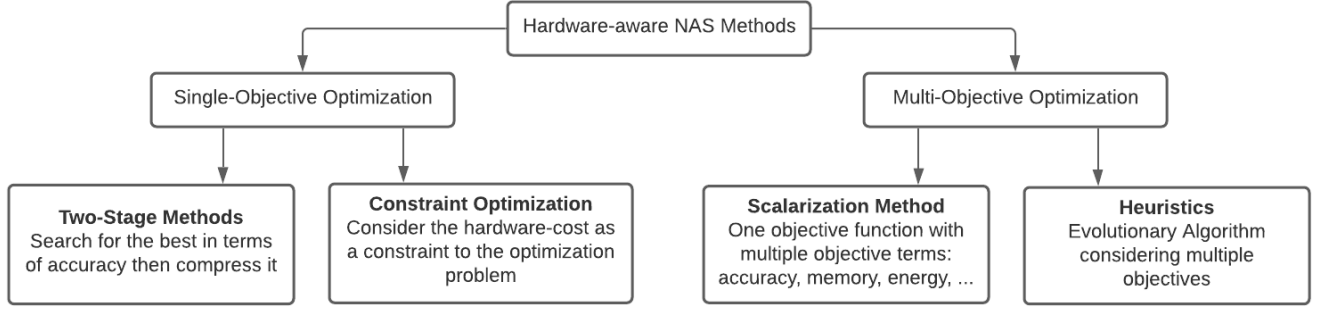


Fig. 3. Problem Formulation of the NAS process [1]

configurations where most hardware-aware NAS belongs to. Single target means that the neural architecture search tries to find the best architecture for one single hardware target. We divide this category further in hardware-aware search strategy and hardware-aware search space. The *hardware-aware search strategy* focuses on solving the neural architecture search as a multi-objective problem which means that the accuracy is taken into account as always but also some hardware measurements for example latency. The opposite to this is the *hardware-aware search space* where the neural architecture search works with only a pool of architectures. The bad working architectures on the target hardware are eliminated. It helps to have a prior knowledge for the target platform for creating the set of architectures. The neural architecture search considers then only the accuracy for this set of architectures and no hardware measurements.

The second category of goals of hardware-aware neural architecture search is **single target, multiple configurations**. This category not only focusing on getting the highest accuracy. The goal is to find the best architecture in terms of the hardware specification in combination with the accuracy. So, one has multiple platform configurations and try to find the best one.

In addition, it exists the category **multiple targets**. This means one does not focus on one specific hardware like FPGAs but search for the best architecture for multiple hardware platforms. This approach is the most desirable one because it allows one to port the best architecture to different hardware platforms. But it is also the most challenging one because architectures perform differently on different hardware. We also consider this problem for having an architecture for multiple platforms in Section X.

VI. TARGET HARDWARE PLATFORMS

There exists three different categories of target hardware platforms. Firstly the **server processors** which can be CPUs, GPUs, FPGAs and ASICs. Its goal is to maximize the accuracy. There is a increasing interest in adapting the neural architecture search for a specific hardware for reducing training costs, but mostly up to now researcher do not take into account the constraints of each hardware platform. The second

category are **mobile devices** which are widely used but do not have much memory and computational power. Therefore, more and more research is focused on gaining high accuracy by using constrained hardware. The last category are **tiny devices** which are devices that run deep learning model with very low power.

VII. SEARCH SPACES

There exist two search spaces used by the search strategies of hardware-aware neural architecture search.

A. Architecture Search Space

The first one is the Architecture Search Space. It defines a set of architectures which means a set of possible operators and their connections. For designing an architecture search space one can use one **fixed architecture** and only search for the best architecture hyper-parameters. Unfortunately this needs some knowledge of the human which designs the architecture, but it reduces the size of the search space. On the other hand one can do a search without a fixed architecture and search for **fitting operators and their connections** of a set of operators and connections. This does not need knowledge of the human how the architecture structure should look like but of course it produces a bigger search space size and with this also a longer search time. By searching for the best architecture without a fixed architecture one distinguishes between three types. Firstly one can have a *layer-wise search space* where the architecture is searched out of a set of operators. Another option is to use a *cell-based search space*. A cell is defined as a graph that includes some operators. A cell-based search space creates a architecture by repeating a founded cell multiple times. This mostly leads to a high accuracy but it limits the flexibility in hardware specialization by stacking the same cells. In addition, there exists the *hierarchical search space*. The architecture is build by constructing blocks out of a defined number of different cells. This solves the missing flexibility of the cell-based search space because one does not only stack the same cells on top of each other.

In the past most NAS researchers defined the type of a network which should be founded by the algorithm. Possible types are convolutional neural networks (CNN's) or recurrent

neural networks (RNN's). The trend now switches to not restrict the search too much in one network type.

B. Hardware Search Space

The second possible search space is the Hardware Search Space. This means before evaluating the architectures one already does some hardware optimizations. Unfortunately, it is not possible to add all hardware specifications in the search space for optimizing because it would explode the search space and in combination the search space time. So, beforehand hardware optimizations cannot consider all hardware specifications.

In general, there exist two different types of hardware search space. The search space can be **parameter-based** which means that the search space is a set of optimizing parameters, e.g. tiling parameters or buffer sizes. On the other hand a search space can be **template-based**. So the search space consists of a set of pre-configured templates instead of optimizing parameters. The idea is to use already known successful designs as templates.

VIII. SEARCH STRATEGIES

The search strategy of hardware-aware NAS consists of three elements: an accuracy evaluation method, a hardware cost evaluation method and a search algorithm.

A. Accuracy Evaluation Method

The accuracy is evaluated by training the architectures of the search space and comparing the accuracies. Because it would be very time and computational consuming to train all architectures completely, there exists some methods for estimate the accuracy without a training. First, one can **share weights** between different models to reduce the search time. Second method for speeding up the determination of the accuracy is **early stopping**. This means one train a deep learning model only a few epochs and use the accuracy as an approximation for a complete training. In addition, one can include **hot start**. Hot starting an architecture search means that the searching and training of an architecture is not started with a random model but with an already known efficient model. Furthermore, **proxy data set** can be used to only train the deep learning model with some data elements and increase the number in the last steps. In the end there also exists already some **accuracy prediction models** which use the architecture and data set for predicting the accuracy beforehand.

B. Hardware Cost Evaluation Method

Hardware cost evaluation methods are methods for measuring hardware metrics in real-time or as an estimation. A possible hardware measurement are the **FLOPs and Model Size**. FLOPs are floating point operations. This means the number of floating point operations and the number of parameters as the model size are used as hardware cost measurement. Unfortunately, the number of FLOPs do not completely correlate with the execution time. This means, FLOPs are not a good indicator to search for efficient architectures,

but a small model size leads to fewer memory consumption and automatically searches for compressed models. So, the combination of FLOPs and the model size can be used as a hardware measurement. Another option is to measure the **latency**. A low latency is important for using NAS in real world applications for ensuring a fast action. Also, the **energy consumption** measurement can be used as metric for hardware cost evaluation. Possibilities are to use the peak of an energy consumption or the average consumption. Furthermore, the **area** means that it is desired to take into account to gain the smallest possible chip by evaluating the hardware cost. The area also correlates with the power consumption as well. In the end the **memory footprint** can be considered as a hardware measurement. As said above the number of parameters from a model is a good indicator for measuring the memory footprint. Another better option is to measure the memory consumption while running.

C. Calculation of FLOPs and Model Size

In our future work we will use the FLOPs and Model size as the numbers of parameters as our hardware cost measurement. Therefore, we present the calculation of FLOPs for convolutional layer and fully connected layers which are the layers we will use in our future work very likely [5]. The calculation of FLOPs for convolutional layer is shown in equation “(3)”.

$$FLOPs = 2HW(C_{in}K^2 + 1)C_{out} \quad (3)$$

In the formulation equation “(3)” stands H for the height of the input, W for the width of the input, C_{in} are the number of input channels, K for the kernel width and C_{out} is the number of output channels. We assume a symmetric kernel width.

The calculation of FLOPs for the fully connected layer is presented in “(4)” whereby I is input dimensionality and O is output dimensionality.

$$FLOPs = 2(I - 1)O \quad (4)$$

On the other side the calculation of the model parameters for convolutional layers is done as follows [6]:

$$Params = C_{in}K^2C_{out} + C_{out} \quad (5)$$

Again, the K is the kernel width which is assumed symmetric. C_{in} is the number of input channels as well as C_{out} the number of output channels. When there is no bias, the formulation “(5)” is reduced to $Params = C_{in}K^2C_{out}$.

The formulation of calculating the number of parameters for a fully connected layer is shown in “(6)” whereby again I is the input dimensionality and O is the output dimensionality.

$$Params = IO + O \quad (6)$$

1) *Real World Measurements versus Estimation models:* Of course real world measurements are very good for evaluating the hardware cost because we are getting a very precise accuracy in comparison to the actual using. But real world measurements have some drawbacks. Mostly, it is very costly. One has to have all used hardware platforms available. In addition, the searching for the best architecture runs longer when trying to get more precise accuracy. Therefore, there exists estimation models like prediction models, lookup tables or a computation of an analytical estimation. In general, prediction models speed up the search time the most. After this follow the lookup tables and then the analytical estimation. But the problem of creating some of these estimation models is that the developer needs hardware experiences and knowledge for making an estimation about the hardware cost.

D. Search Algorithm

The search algorithm samples the architectures from the search space and updates the search space for getting architectures with a higher accuracy. The two most used search algorithms in neural architecture search are reinforcement learning and evolutionary algorithms. By using **reinforcement learning** an agent chooses between actions in an environment and gets a reward after each action. Its goal is to maximize the reward so the decision of next actions are based on this. By using reinforcement learning for hardware-aware NAS one action is the sampling of one architecture and the reward is the accuracy and the hardware cost. Therefore, the sampling depends on gaining a high accuracy and low hardware costs. In contrast to this are the **evolutionary algorithms**. They have three main characteristics: population-based, fitness-oriented and generations. So in the first step the population-based characteristic is applied, which means that the algorithm starts with a set of candidate architectures which are known as population. Then fitness-oriented means that each solution of the set has a quality which is expressed with its fitness score. For improving this fitness score for different architectures, the evolutionary algorithms can make mutations and crossover operations for generating a new generation/population. The best architectures are used as the new starting point for making new populations and calculate again the fitness score.

Other possible search algorithms are for example **gradient-based methods**. Instead of separating the search and the evaluation one creates a super-network which can simulate any child model. So parts of the model can share their weights which reduces the search time. By training the super-network one gets the weights and parameters for all child architectures. Another option is **random search** which is the easiest way by just randomly searching an architecture but it consumes a lot of time in searching.

E. Non-differentiable Hardware Constraints

By using super-networks like in the gradient-based methods the loss function and the hardware cost function must be differentiable. So for getting a differentiable architecture

search space one can use for example Gumbel Softmax, Estimated Continuous Function or REINFORCE algorithm.

IX. NAS FRAMEWORKS

There are multiple frameworks for doing neural architecture search, for example Auto-Keras, Google AutoML, IBM NeuNets, TPOT, Microsoft Archai, deci.ai AutoNAC, Darwin and Microsoft NNI. All these frameworks have their advantages and disadvantages and focus on different use-cases. We focus on the **Microsoft Neural Network Intelligence (NNI)** framework because it seems to fit the best for our future work. Microsoft NNI can be used to better compare, reproduce and experiment with various NAS algorithms. It helps the researcher to design their neural network architecture for example by helping with the definition of a super network. It concentrates on the efficiency of the automation by doing NAS. In addition, researchers can use existing NAS algorithms and can easily modify them. Except for AutoNAC, all existing frameworks are not hardware-aware, but Google AutoML and Microsoft NNI support compression methods for specializing the target platform.

X. CHALLENGES AND LIMITATIONS

In the following we list various challenges and limitations of hardware-aware neural architecture search.

A. Benchmarking and Reproducibility

For providing comparison, it should be possible to reproduce a NAS process. But this is difficult because of the huge possibilities of search spaces, training methods, computational resources and especially for hardware-aware NAS possibilities of hardware devices. For solving this problem there are set some benchmarks. Firstly, for reducing the needed costs researchers should directly query a tabular data set instead of generating a search space. In addition, it should be used various search strategies on the same search space for a reliable comparing between the strategies. Furthermore, data sets which are used by accuracy predictor models and hardware cost models should be provided. Also, make NAS easier for developers without hardware knowledge by proposing data sets with hardware-related metrics.

B. Transferability of AI Models

Transferability here means to train a model on one data set and fine-tune it after that on the target data set. Reasons for doing this are that one can use a simpler and smaller train data set or if the target data set has not enough data points for training.

Most researchers use cell-based search space in neural architecture search for transferring their model. For transferring, you modify the network by adding more cells of your defined cells and adjust the input sizes of the cells. Because of bad results by stacking the same cell on top of each other, some recent works start to transfer hierarchical search space by using different operators and combinations for the cells. Another option is to discard the fine-tuning and directly find

transferable weights. For this one create a task-specific super-network from an overall super-network and search then for architectures within this network. Therefore, no extra training is needed.

In general, transfer a model in NAS is very difficult. A starting point is to create a macro-architecture for a specific task and modify the search space for this task. By considering the hardware, one has to transfer to different hardware platforms which is on top very difficult as described in the next Section.

C. Transferability of the hardware-aware NAS across multiple Platforms

One interesting topic in hardware-aware neural architecture search is to design an architecture and transfer it to multiple different platforms. The problem hereby is the variety and different complexities of the hardware platforms. For solving this issue there exists two approaches.

1) *Transfer the entire NAS Process:* The idea is to transfer a NAS to a different hardware platform by changing the measurement values. For this the whole NAS process is run again on the different hardware. This rerun of the NAS process leads to high computational consumption's. In addition, it can be a huge effort to collect the hardware constraints. The hardware measurements can be done with real-world measurements, analytical estimation, lookup tables and prediction models like explained before. But real-world measurements need the target hardware platform already during search time and are very slow. Analytical estimation needs experiences from the developer in the specific target hardware whereas lookup tables and prediction models have to run the complete set of operators again. Therefore, all in all it is very difficult and not scalable to transfer a NAS process to a new hardware by changing the measurement values.

2) *Transfer the final Model:* The second approach is to transform a final model of a NAS process to fit to a different hardware. For fitting the final architecture to a smaller device, mostly compressing the model is used. But this specialization has some limits. Firstly, different operators are not equally efficient for two different platforms. This also makes the subgroup of hardware-aware NAS which tries to find the best architecture for multiple hardware platforms difficult. In addition, there are limits in the compression methods. Compressed models never lead to a higher accuracy as the original model.

XI. CONCLUSION

This paper gives a brief overview about hardware-aware NAS for quantized networks. We presented the definition of NAS and the possibility to define hardware-aware NAS as a single- or multi-objective optimization. Search Strategies and Search Spaces are discussed as well as various hardware platforms and NAS frameworks. In the end we consider existing challenges and limitations of hardware-aware NAS today.

A. Future Work

The idea of our future work is to do hardware-aware NAS for quantized networks. This means our model optimization techniques are focusing on quantization. Maybe we will also add pruning for making the deep learning model smaller. The idea is to use hardware-aware NAS with these model compression methods for designing the model. There already exist some researches which

We did not decide until now if the NAS process is done as a single-objective optimization or a multi-objective optimization. In general it would be easier to optimize a single objective, i.e. the accuracy because we then mainly focus on the accuracy as the only objective and take the hardware usability after the designing or only with some threshold into account. But a multi-objective optimization would be more precise because it can be that single-objective optimization do not perform well on the specific hardware. So a multi-objective optimization would be desired.

Our goal is to not focus on a single target, therefore we use the FLOPs and the model size as our hardware cost measurement. But it is not really solved until now to find universal deep learning models which perform very well on different hardware platforms. We discussed the limits of transferring hardware-aware NAS across multiple platforms in Section X. So if we had to decide a hardware platform for gaining reasonable optimizations, we will choose FGPA's. It is not clear until now if we consider a fixed platform configuration or multiple platform configurations for our target. In addition, we focus on server processors where FPGA's belong to and not mobile devices or tiny devices.

For clarifying which search space and search algorithm we want to use, we need more research in recent works which use hardware-aware NAS to determine the best search space and search algorithm for our use-case. For now we have a lack of knowledge about different algorithms and search spaces. We are open which accuracy evaluation methods we will use for estimating the accuracy without long training. To measure the hardware cost we will take the number of bit operations or floating point operations and the number of parameters as model size into account. We are not able to real world measurements, if we stay independent of the hardware platform. Therefore, we have to use estimation models for measuring the hardware cost.

As already mentioned our used NAS framework is Microsoft NNI¹. In addition our machine learning framework will be PyTorch [4]. Unfortunately there is currently no official support for calculating FLOPs with PyTorch², so we have to calculate it by our own or with third-party libraries. Gaining the model parameters is provided by PyTorch.

Missing: verwandte Arbeiten, distiller, co2

¹<https://nni.readthedocs.io/en/stable/index.html>

²<https://github.com/pytorch/pytorch/issues/5013>

REFERENCES

- [1] Hadjer Benmeziane, Kaoutar El Maghraoui, Hamza Ouarnoughi, Smail Niar, Martin Wistuba, Naigang Wang, "A Comprehensive Survey on Hardware-Aware Neural Architecture Search", 2021.
- [2] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012): 1097-1105.
- [3] Collobert, Ronan, et al. "Natural language processing (almost) from scratch." *Journal of machine learning research* 12.ARTICLE (2011): 2493-2537.
- [4] Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. "PyTorch: An Imperative Style, High-Performance Deep Learning Library." In: *Advances in Neural Information Processing Systems* 32 [Internet]. Curran Associates, Inc.; 2019. p. 802435. Available from: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [5] Molchanov, Pavlo, et al. "Pruning convolutional neural networks for resource efficient inference." *arXiv preprint arXiv:1611.06440* (2016).
- [6] Borji, Ali. "Enhancing sensor resolution improves CNN accuracy given the same number of parameters or FLOPS." *arXiv preprint arXiv:2103.05251* (2021).

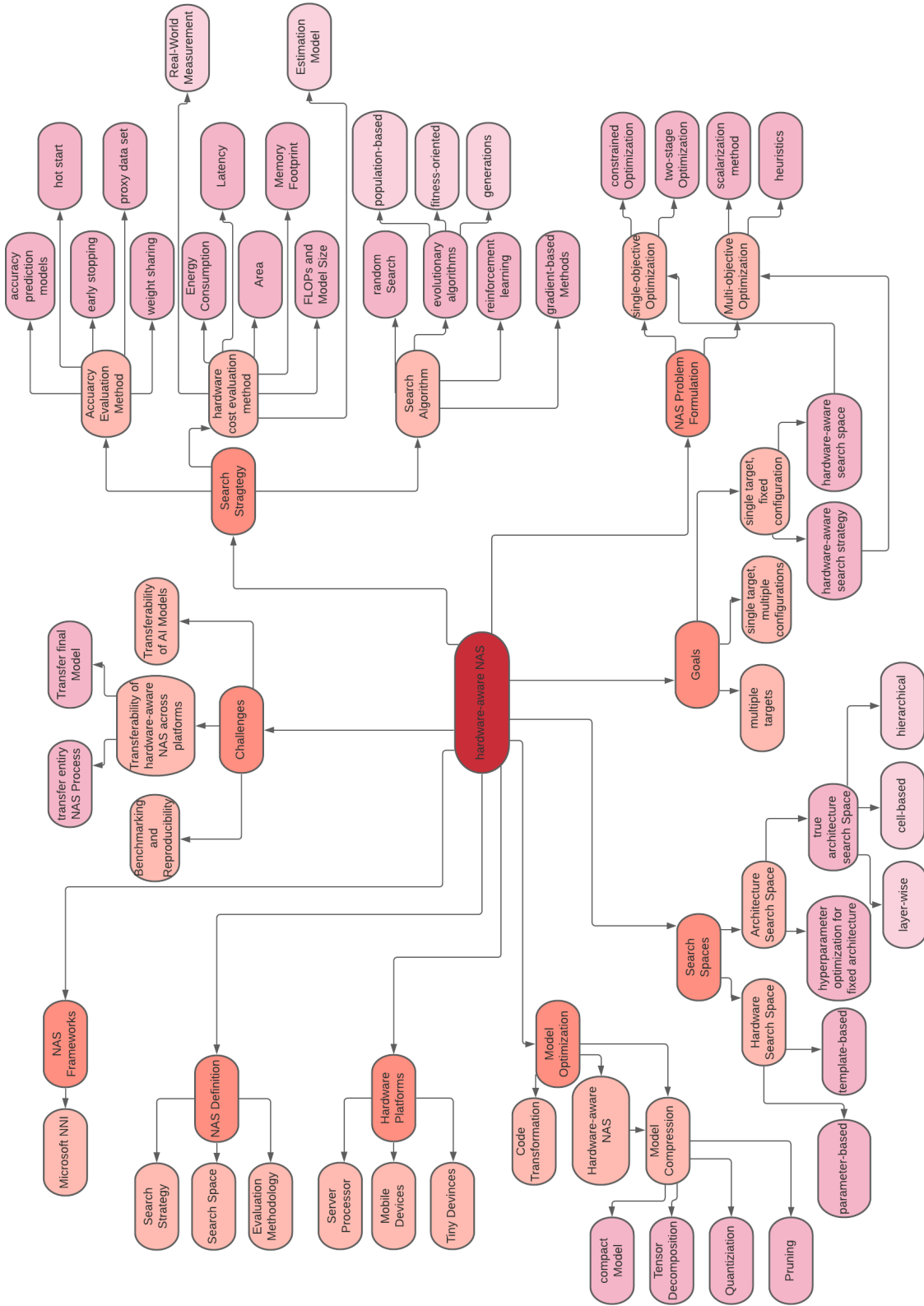


Fig. 4. Mindmap of content of paper