

# Neural architecture search for in-memory computing-based deep learning accelerators

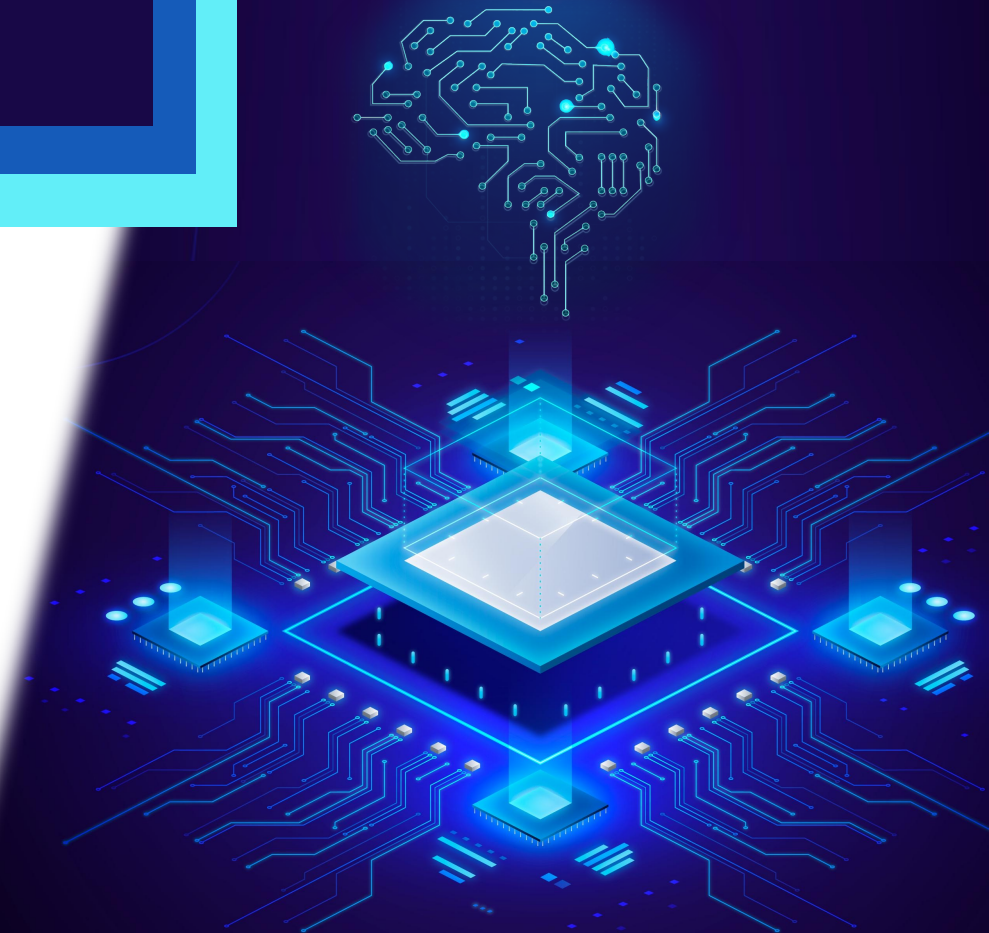
**Olga Krestinskaya**

*King Abdullah University of Science and Technology, 23955, Saudi Arabia*

## **Based on:**

Olga Krestinskaya, Mohammed E. Fouda, Hadjer Benmeziane, Kaoutar El Maghraoui, Abu Sebastian, Wei D. Lu, Mario Lanza, Hai Li, Fadi Kurdahi, Suhaib A. Fahmy, Ahmed Eltawil, and Khaled N. Salama.

*Nature Reviews Electrical Engineering (2024): 1-17.*



# Agenda



## Introduction, motivation and software-hardware co-design

*Motivation to improve hardware efficiency and challenges in selecting optimum design*



## Hardware-aware Neural Architecture Search (HW-NAS)

*HW-NAS methods, algorithms, hardware cost estimation methods, HW-NAS frameworks for IMC*



## Future directions, open challenges and final thoughts

*Roadmap of HW-NAS for IMC, open issues, HW-NAS and other optimization techniques, summary*

# Rapid AI development and Increasing Neural Network Complexity

“AI everywhere”



Vision and language processing



Self-driving vehicles

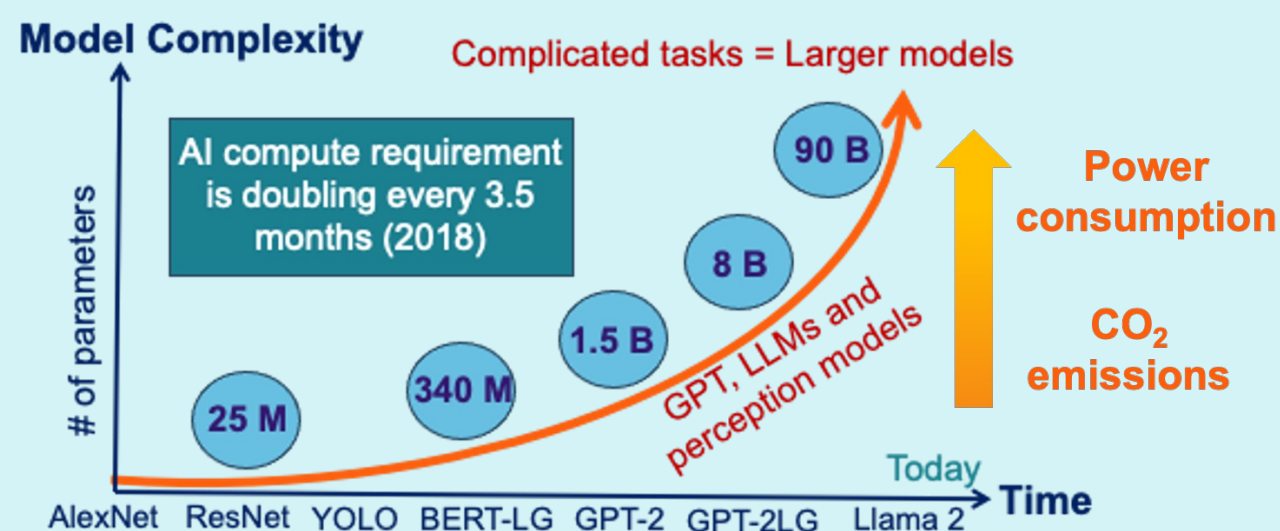


Security applications

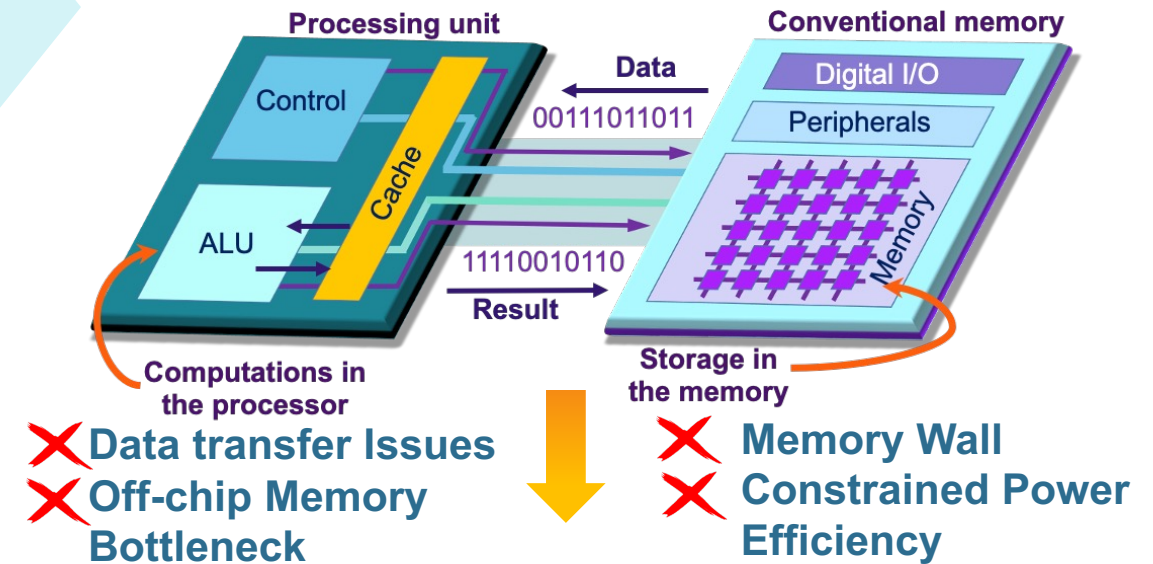


Smart cities

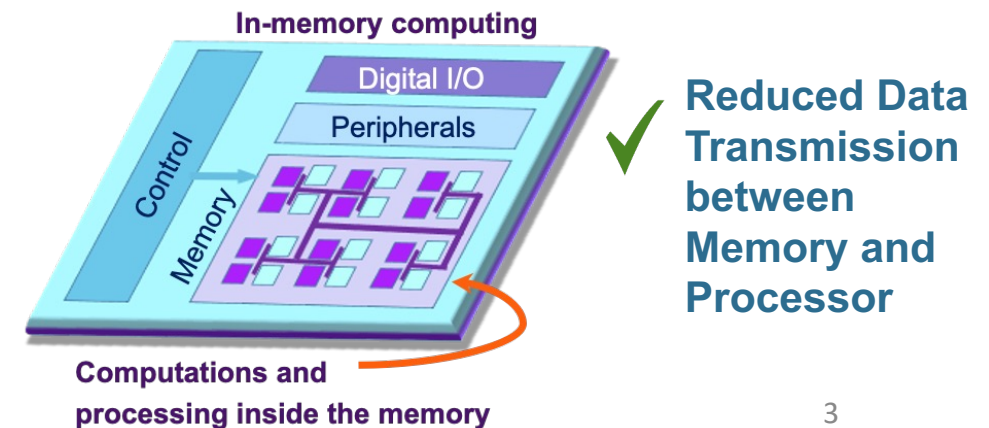
Increasing complexity of neural networks



Traditional von Neumann Architecture:

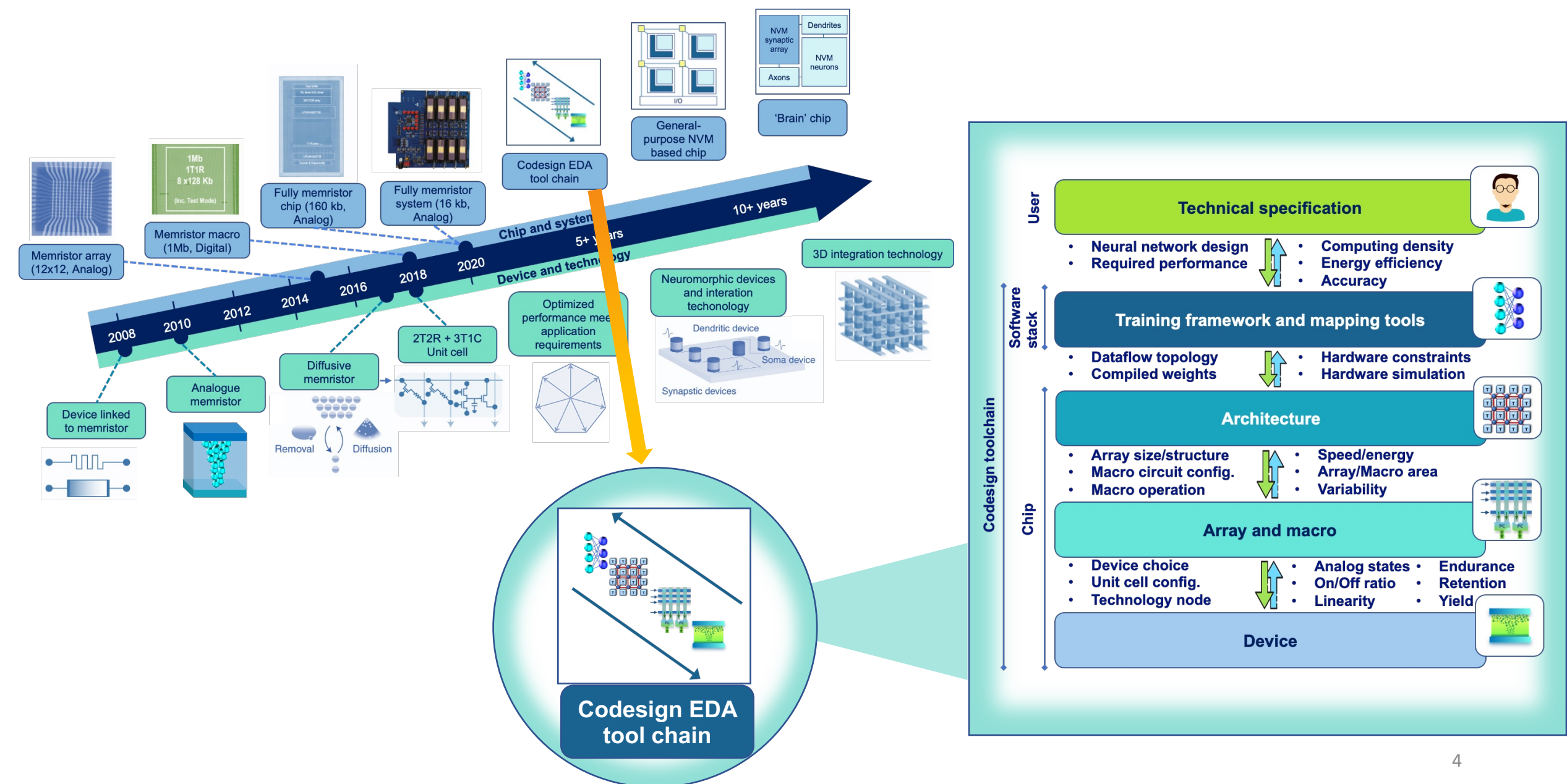


In-memory computing accelerator:



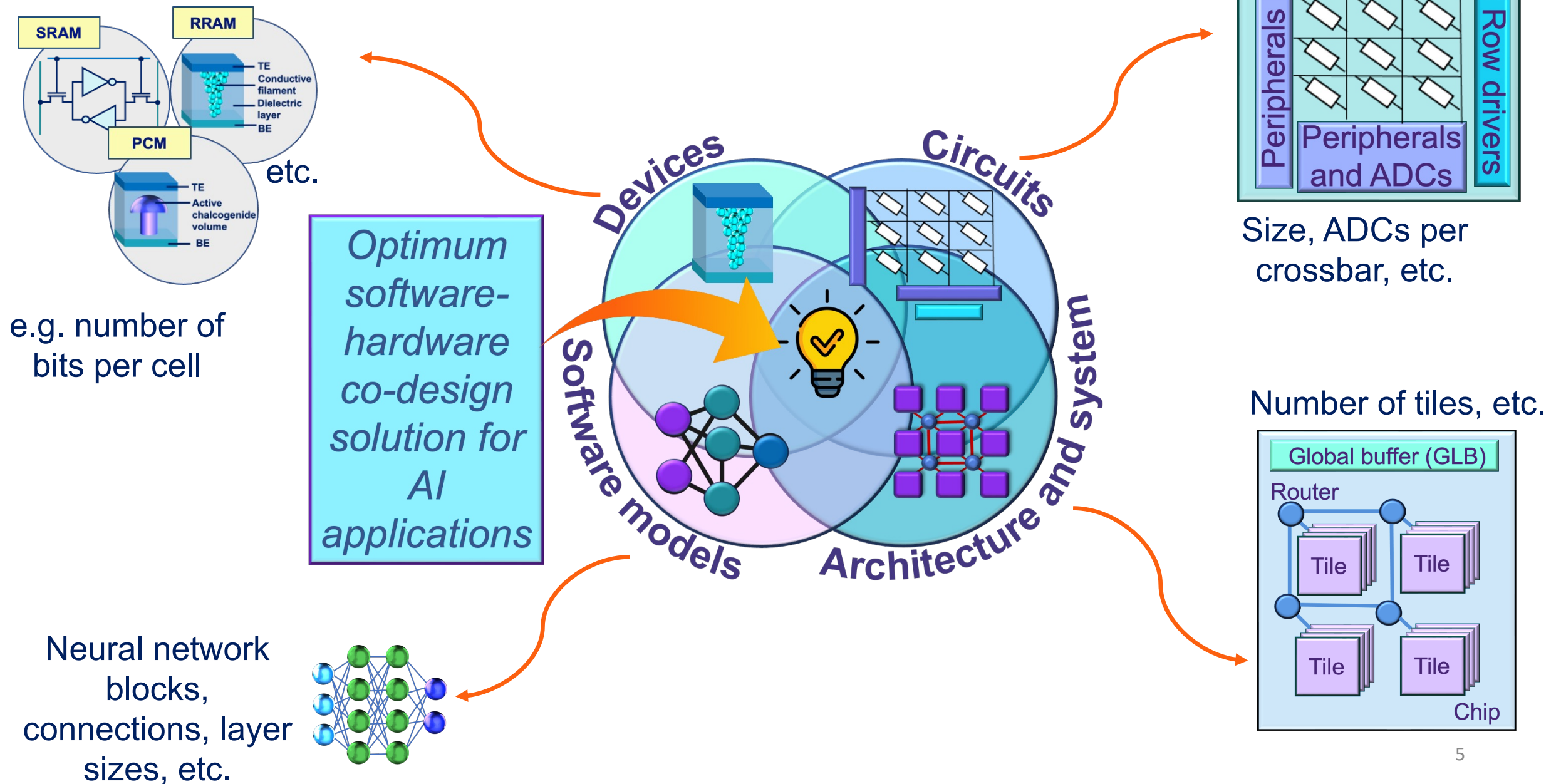


# In-memory Computing Roadmap and Software-Hardware Co-design





# Software-Hardware Co-design for IMC Accelerators

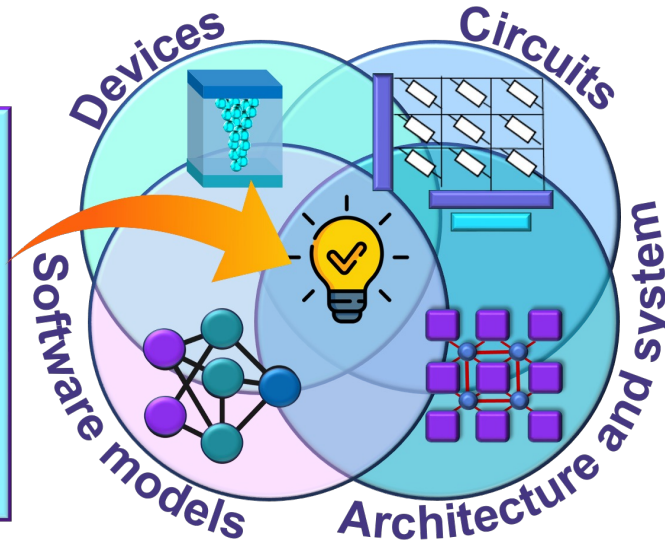


# How to Select the Optimum Design?

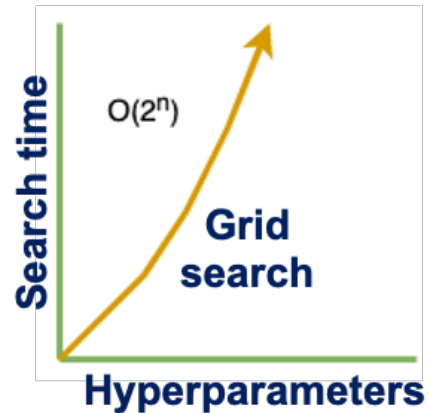


Impossible to  
optimize  
manually

Optimum  
software-  
hardware  
co-design  
solution for  
AI  
applications



One approach to  
achieving such  
optimization is  
through Hardware-  
aware Neural  
Architecture  
Search (HW-NAS).



Too many parameters to consider



Manual optimization of parameters is infeasible,  
based on guessing (sometimes certain rules), and  
require a lot of human efforts



Grid search is slow and search time exponentially  
increases with number of hyperparameters to  
optimize

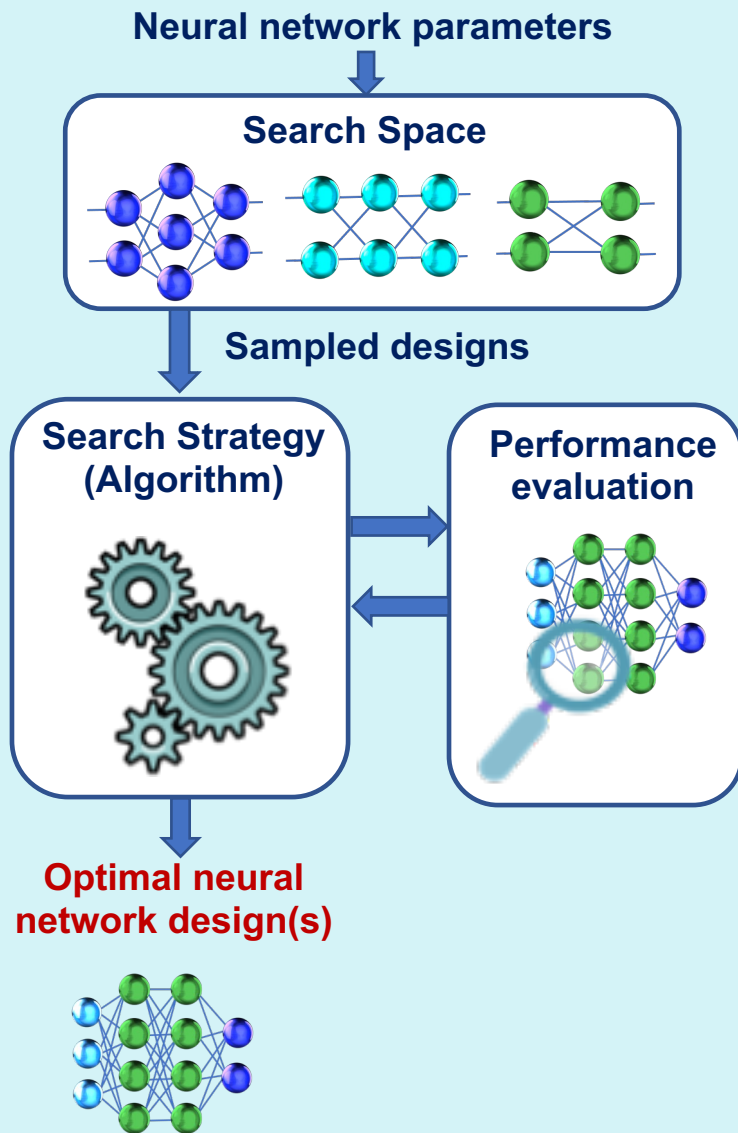


Increasing complexity and network size

Example:  
 $8.5 * 10^{85}$  possible  
combinations

# Neural Architecture Search from Software Perspective

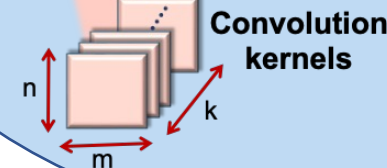
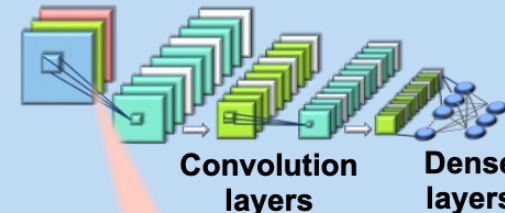
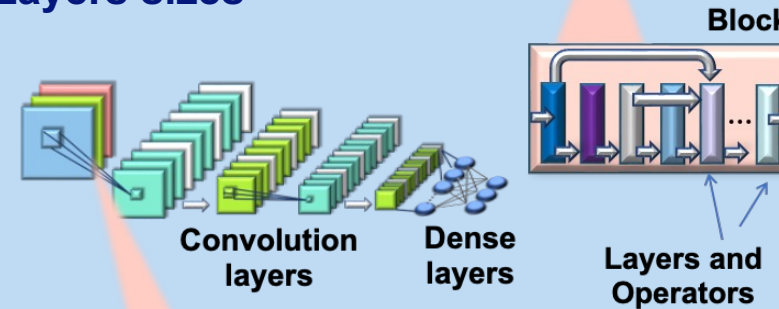
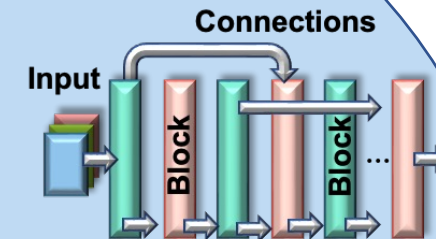
## Neural Architecture Search (NAS)



Often combined with other optimization techniques, e.g. pruning or quantization

## Neural Network Model Search Space

- Layers
- Blocks inside layers
- Kernel size
- Number of kernels
- Layers sizes



Performance evaluation is based on neural network accuracy or loss

Output: neural network(s) of highest performance accuracy



No consideration of hardware efficiency or considering only high-level metrics, e.g. FLOPs



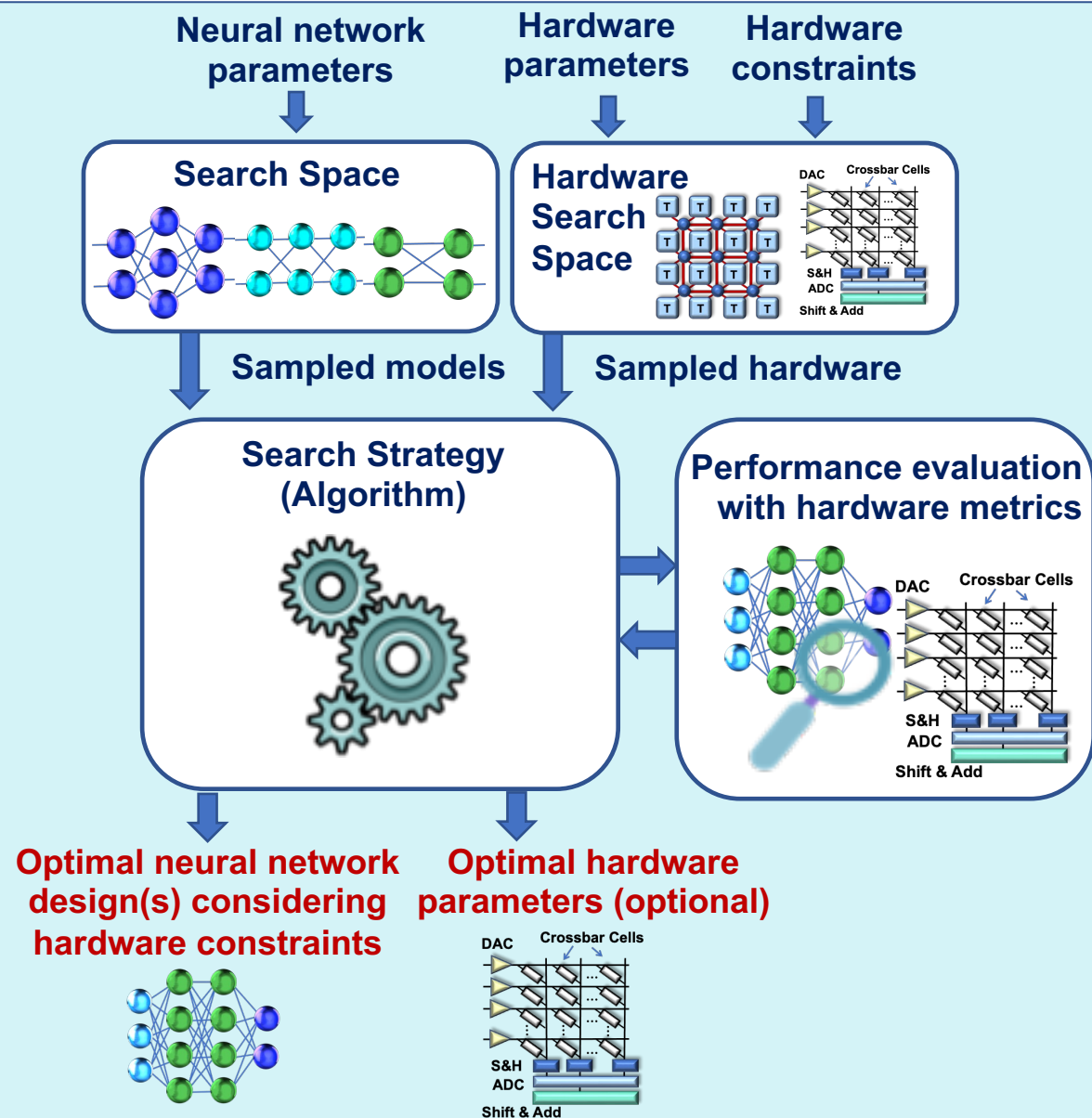
# Hardware-aware Neural Architecture Search

## Neural Architecture Search (NAS)

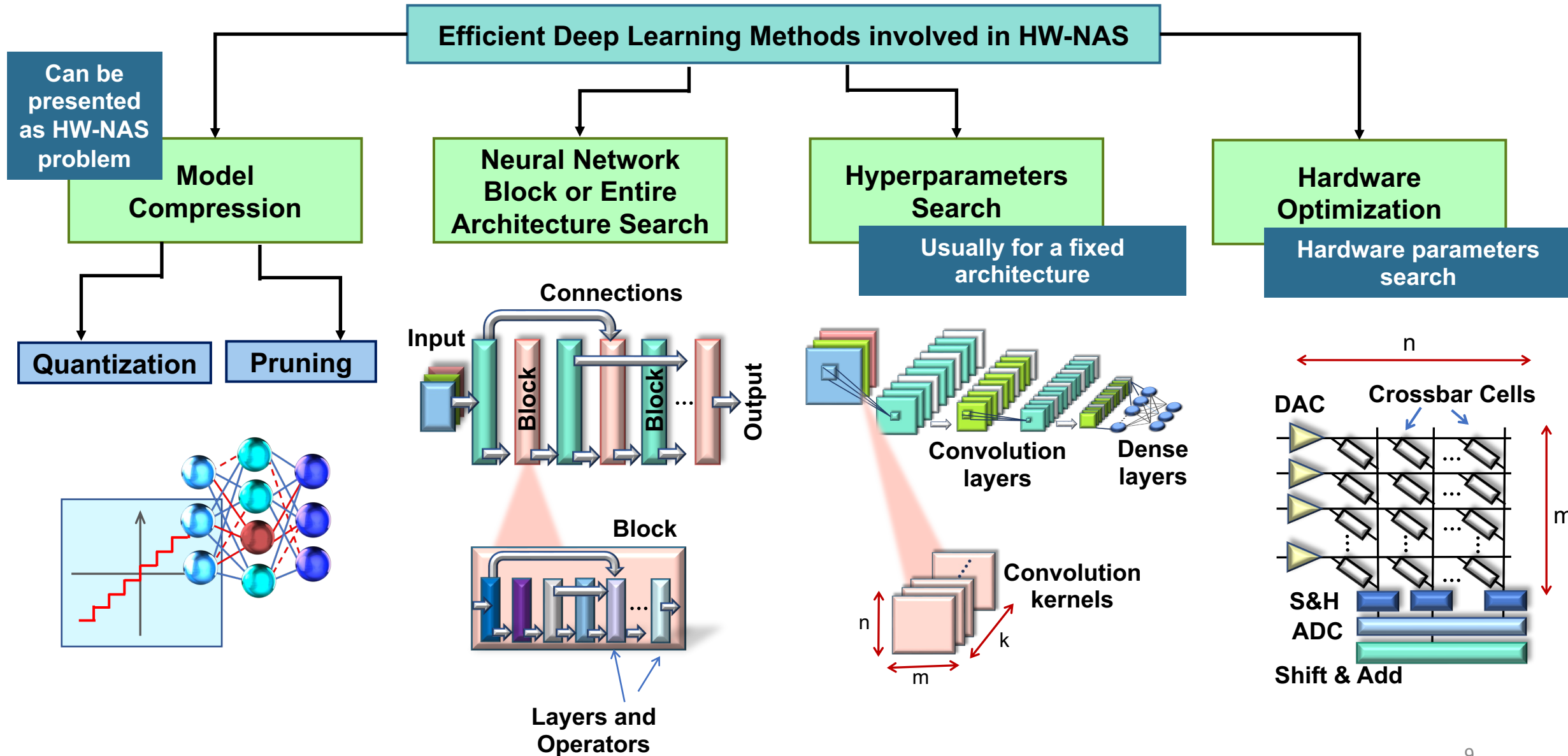


Expanding NAS  
concept to the  
hardware

## Hardware Aware Neural Architecture Search (HW-NAS)

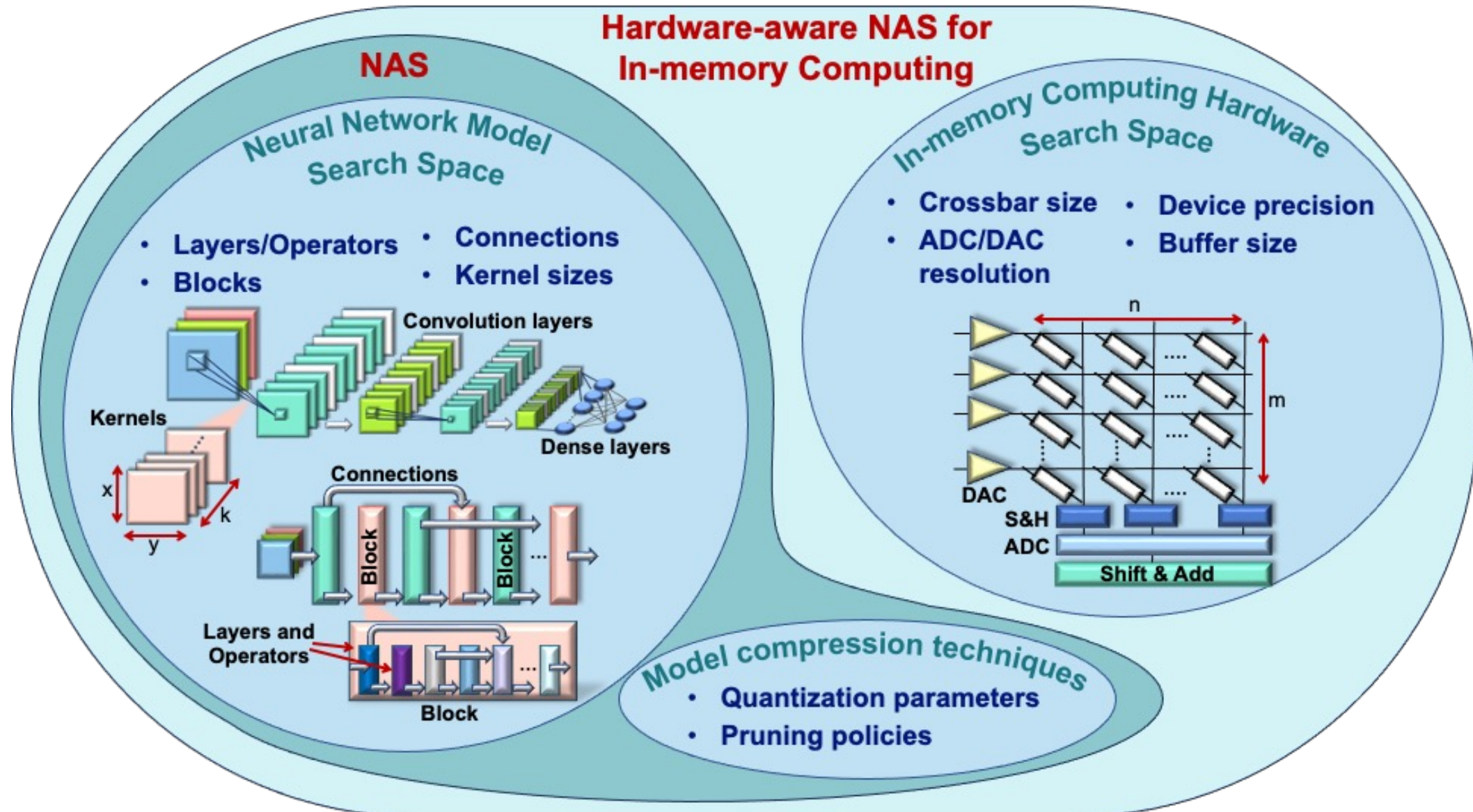


# Hardware-Aware Neural Architecture Search



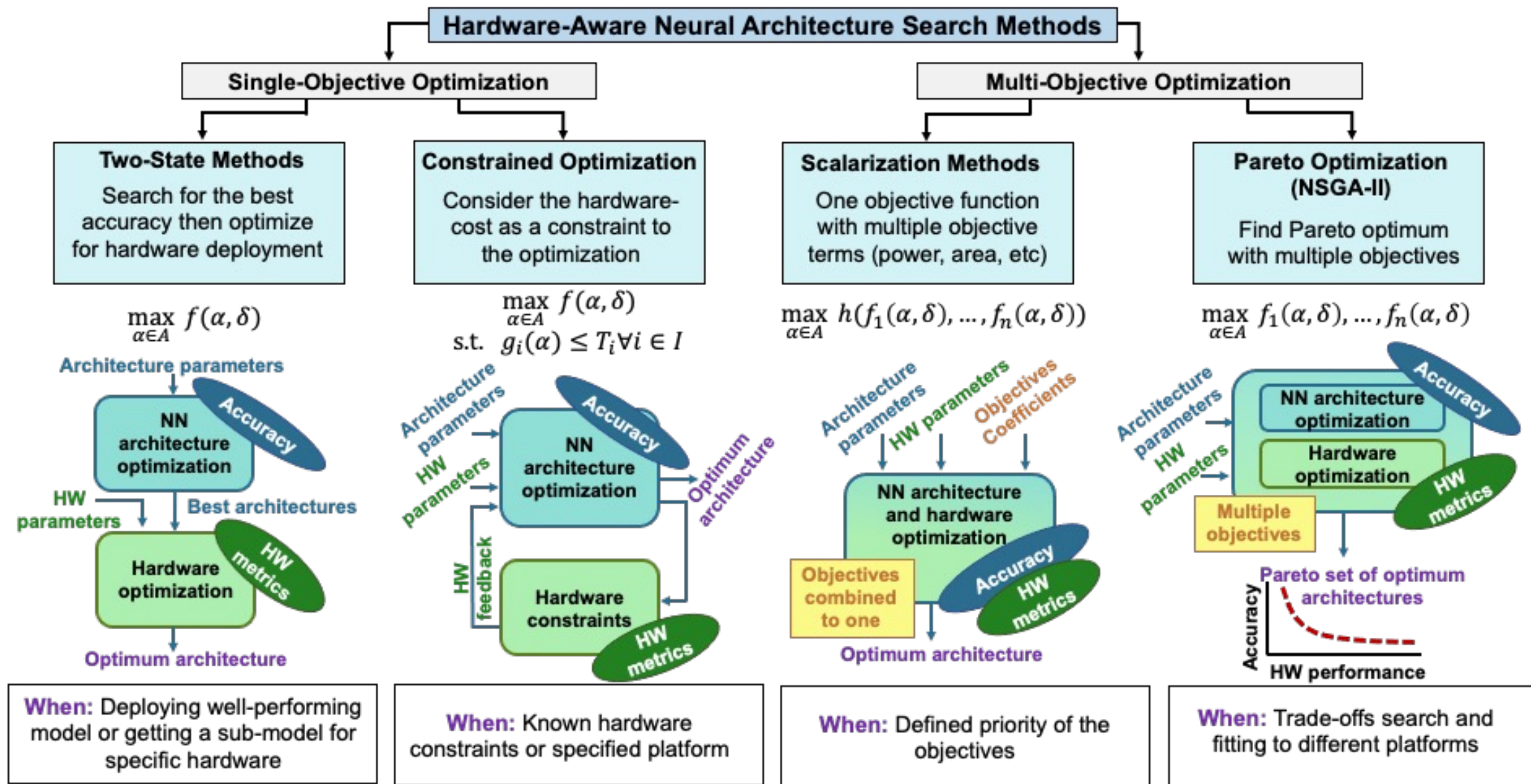
# HW-NAS for In-Memory Computing: Search Space

## Expanding HW-NAS Search Space for IMC Architectures

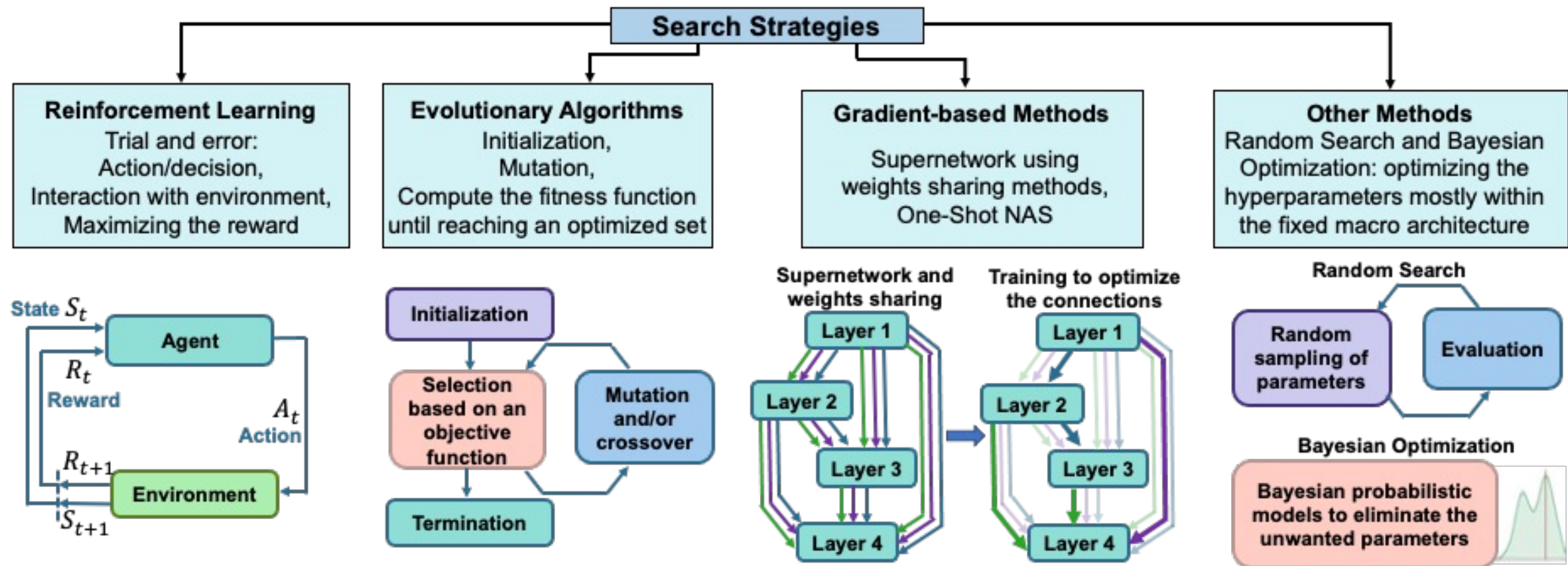




# HW-NAS Methods



# Search algorithms

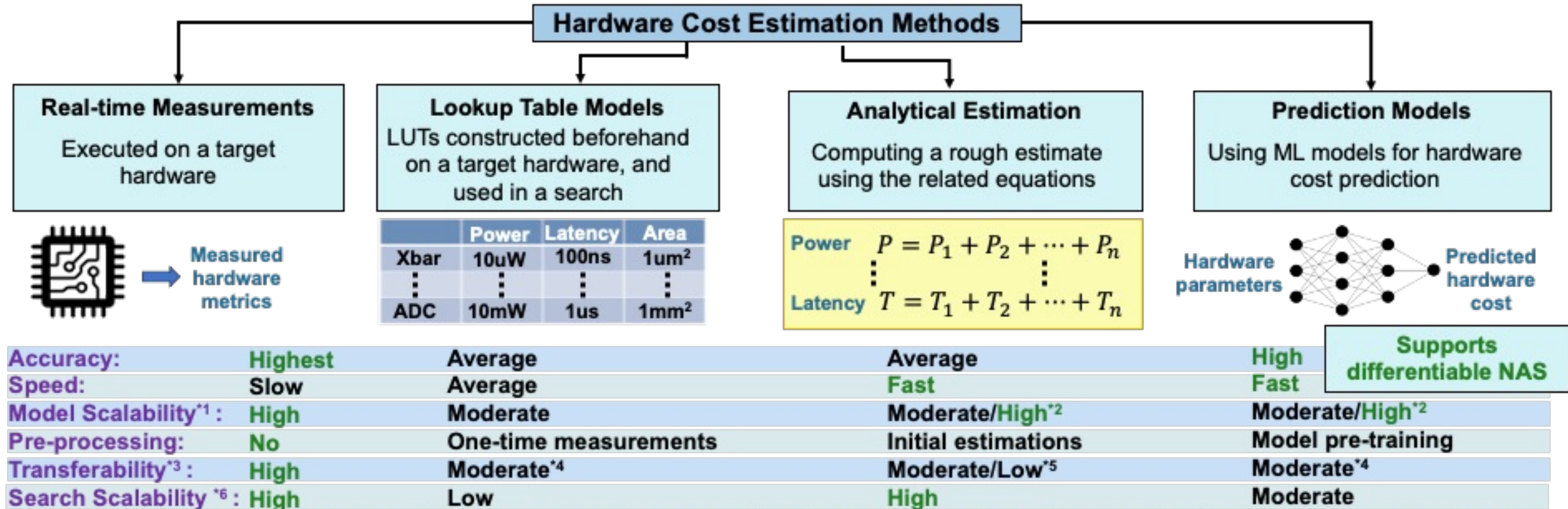


Search Space:	Discrete	Discrete	Differentiable	Discrete
Speed:	Slow	Average	Fast	Slow
Computational demand:	High	Average	Low	High
Memory requirements:	Small	Small	Large	Small
Scalability <sup>*1</sup> :	Not Scalable	Not Scalable <sup>*2</sup>	Scalable <sup>*3</sup>	Not Scalable <sup>*2</sup>

<sup>\*1</sup>: refers to time complexity increasing with a search space, <sup>\*2</sup>: can be used with a supernetwork search space, <sup>\*3</sup>: search time does not increase exponentially with a search space size



# Hardware cost estimation methods



<sup>\*1</sup>: across neural network models, <sup>\*2</sup>: depends on how similar is a new model, <sup>\*3</sup>: across different hardware platforms, <sup>\*4</sup>: requires regeneration, <sup>\*5</sup>: depends on the hardware similarity, <sup>\*6</sup>: with increasing search space (number of hyperparameters in a search)



# State-of-the-art HW-NAS Frameworks for IMC

	Quantization (search)	Pruning	HW- NAS	Architecture Search Space	Hardware Search Space	Hardware cost	Algorithm	Hardware non-idealities
AnalogNAS (2023)	✗	✗	✓	# of blocks, channels, branches, kernel size	✗	AIHWKit	EA	Variations, Cond. drift
NAS4RRAM (2021)	✗	✗	✓	Layers, channels (residual blocks)	✗	RRAM simulator	EA	Variations
FLASH (2021)	✗	✗	✓	# of skip connections, cells, layers, channels	✗	NeuroSim, BookSim	SHGO	✗
NAX (2021)	✗	✗	✓	Kernel size	Crossbar size	GENIEx	DS	Wire/source/sink resistances
Gibbon (2022)	✓	✗	✓	# of blocks, channels, groups, kernel size, bit-width	Crossbar size, ADC/DAC/device precision	MNSIM	EA	Variations
NACIM (2020)	✓	✗	✓	Architecture hyperparameters, bit-width (int./frac.)	Tile/buffer size, bandwidth	NeuroSim	RL	Variations
UAE (2021)	✓	✗	✓	# of channels, filter size, bit-width (int./frac.)	✗	Analytical	RL	Variations, program. errors
CMQ (2022)	✓	✗	✗	Quantization threshold, bit-width	✗	MINT	DS	Variations
Mixed-precision quantization (2021)	✓	✗	✗	Weight/inputs bit-width (int./frac.)	ADC precision	PUMAsim	RL	✗
EGQ (2021)	✓	✗	✗	Weight/activation bit-width	✗	NeuroSim	GA	✗
RaQU (2021)	✓	✗	✗	Weight/kernel bit-width	✗	Analytical	RL	✗
ASBP (2021)	✗	✓	✗	Bits of weights	✗	Analytical	RL	✗
Auto-prune (2021)	✗	✓	✗	Weights (pruned unimportant columns)	✗	MNSIM	RL	✗

⚠ Limited architecture search space (particular network type)

⚠ Limited applications

⚠ Limited hardware search space (only few hardware parameters)

⚠ Consideration of hardware non-idealities and mitigation techniques

⚠ Lack of a unified framework

## The latest frameworks (*not included in the paper but worth mentioning*):

XPert (2023)	Two-step co-optimization of software and hardware parameters, including channel depth, ADC precision, input precision, etc.
CoMN (2024)	Design space exploration for a large IMC hardware search space, including circuits and architecture parameters
Joint Hardware-Workload Co-optimization (2024)	Design space exploration for a large IMC hardware search space to optimize IMC hardware for different workloads simultaneously

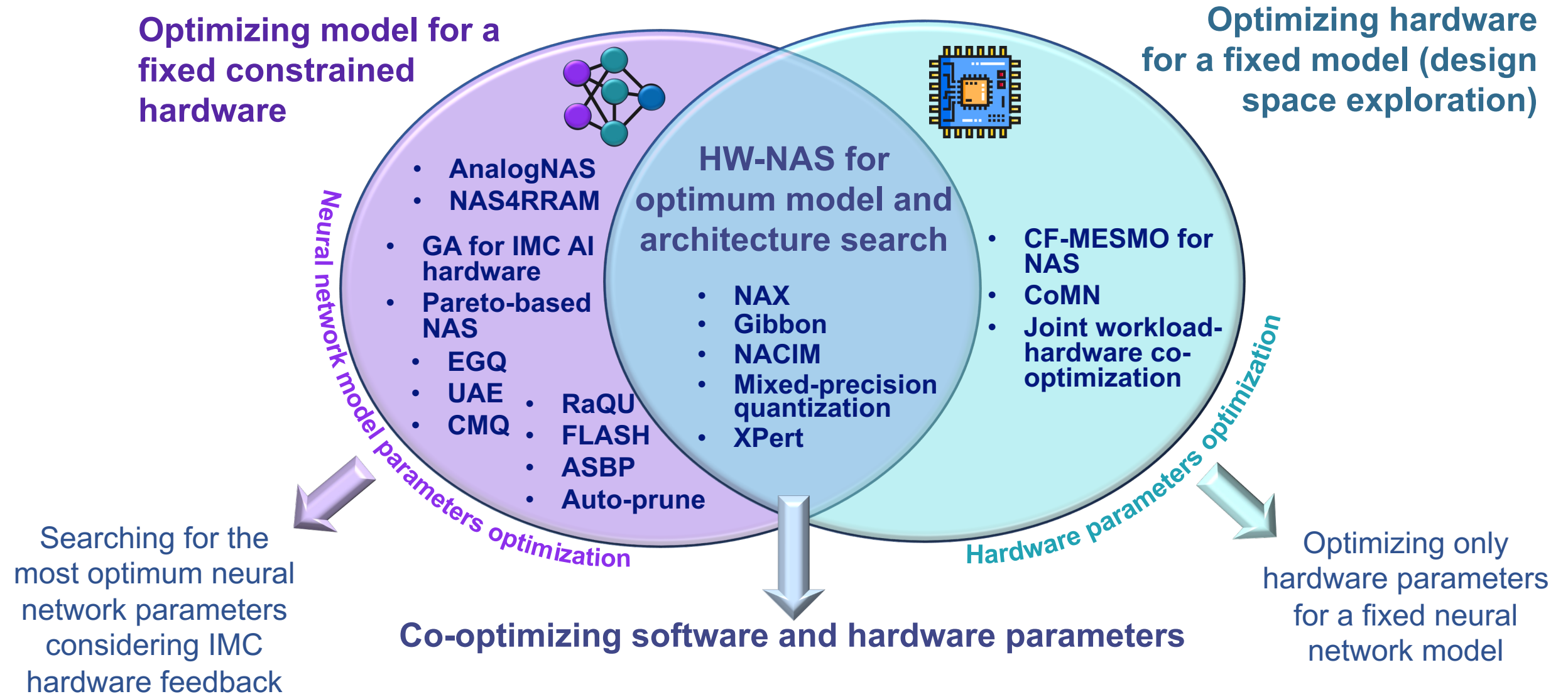
Moitra, A., Bhattacharjee, et al. (2023, July). XPert: Peripheral Circuit & Neural Architecture Co-search for Area and Energy-efficient Xbar-based Computing. In *2023 60th ACM/IEEE Design Automation Conference (DAC)* (pp. 1-6). IEEE.

Han, L., Pan, et al. (2024). CoMN: Algorithm-Hardware Co-Design Platform for Non-Volatile Memory Based Convolutional Neural Network Accelerators. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*.

Krestinskaya, O., Fouda, M. E., Eltawil, A., & Salama, K. N. (2024). Towards Efficient IMC Accelerator Design Through Joint Hardware-Workload Co-optimization. *arXiv preprint arXiv:2410.16759*.

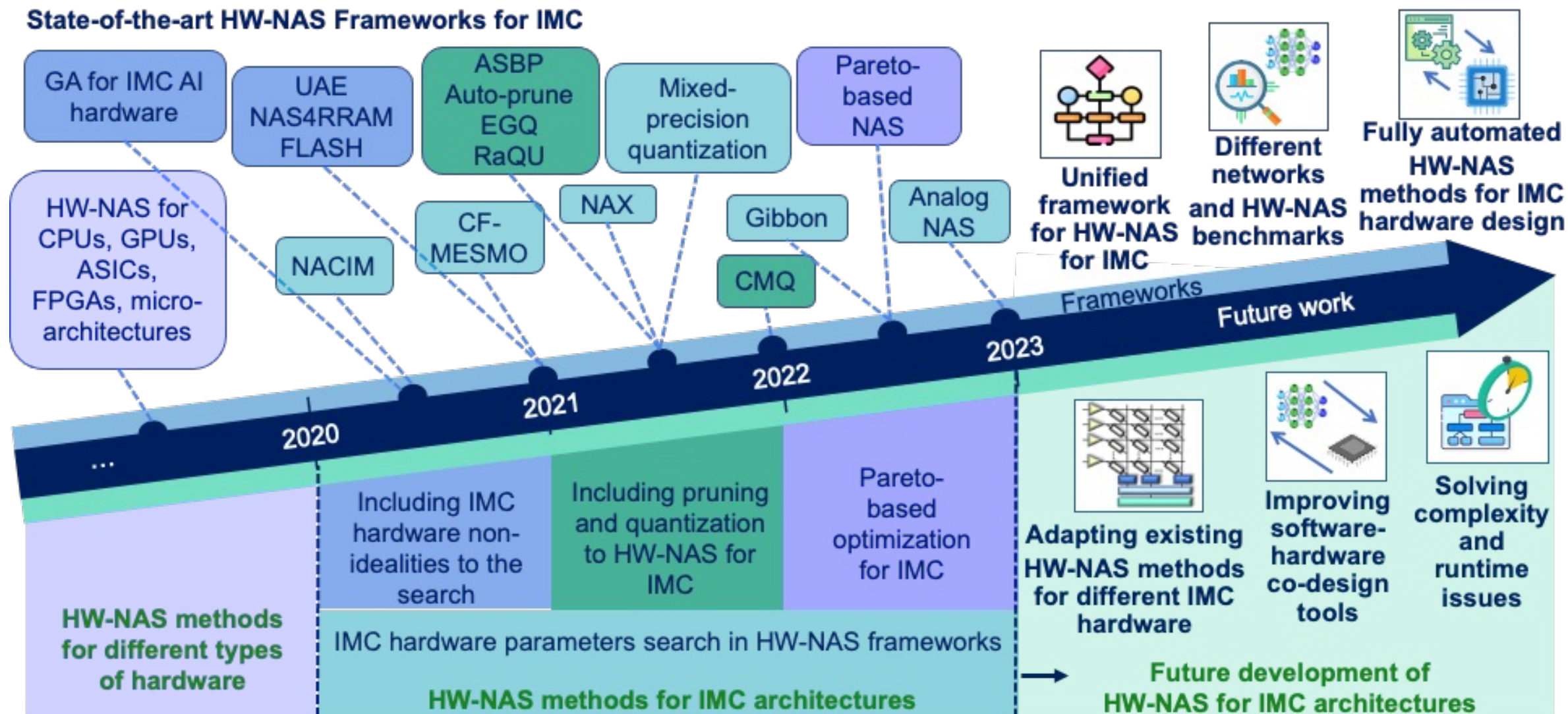
14

# Taxonomy of HW-NAS frameworks for IMC applications



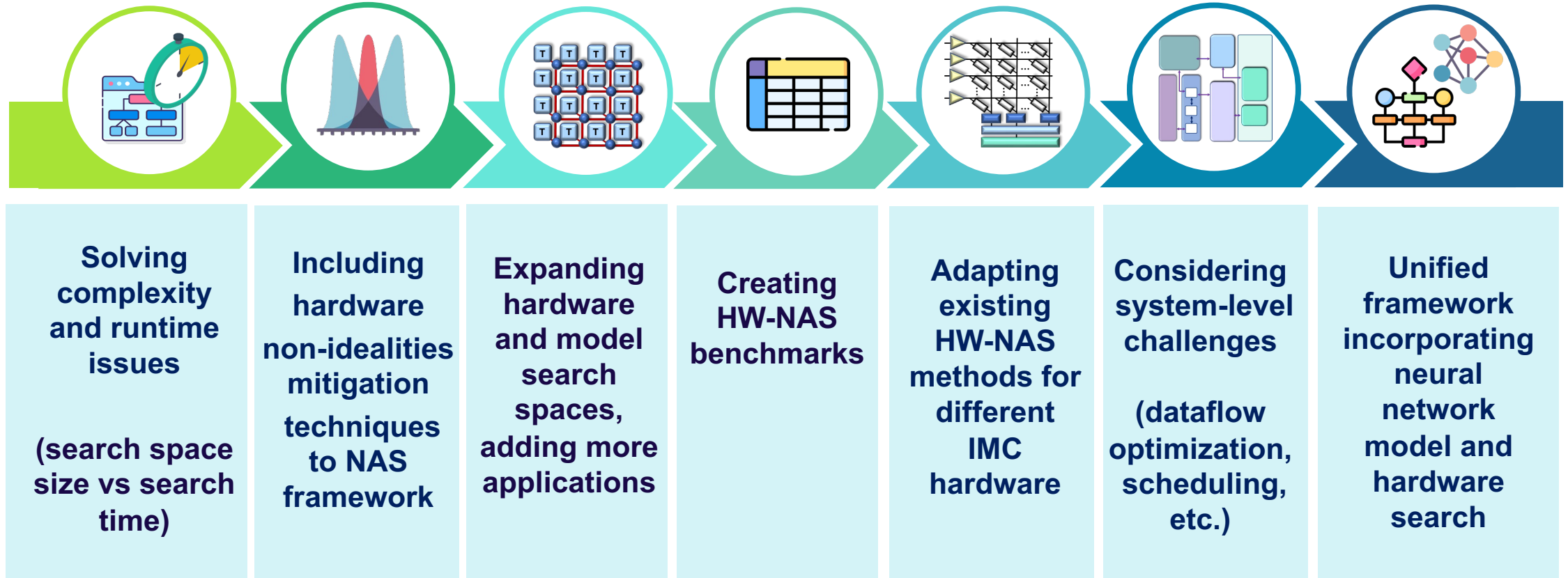
# HW-NAS for IMC Applications Roadmap

## State-of-the-art HW-NAS Frameworks for IMC



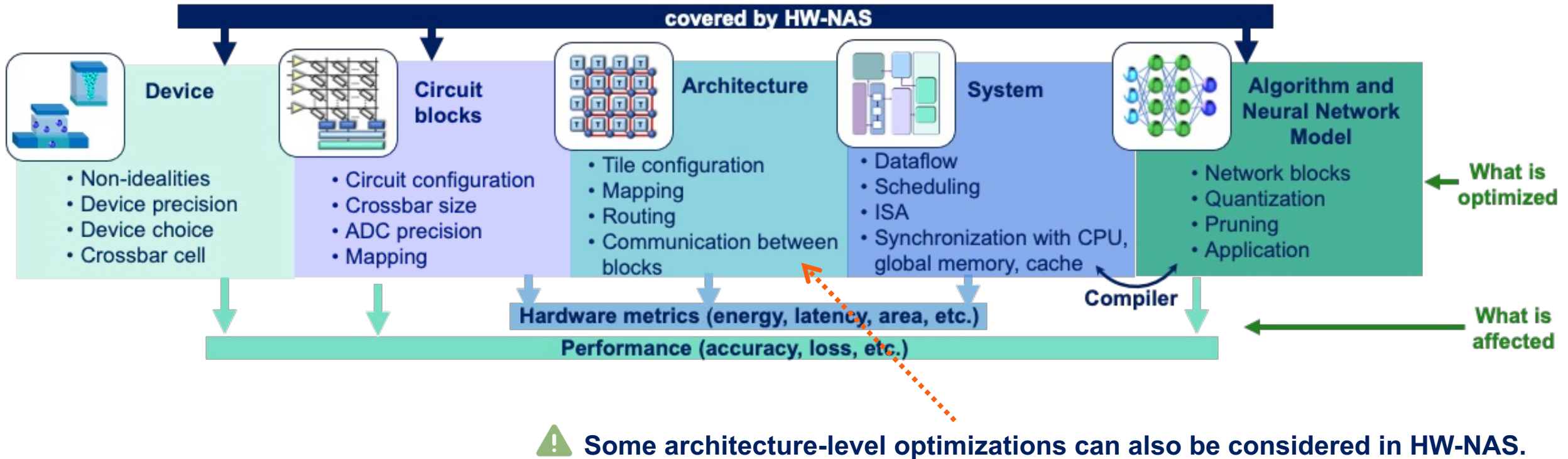


# Open Problems and Way Forward



# Open Problems and Way Forward

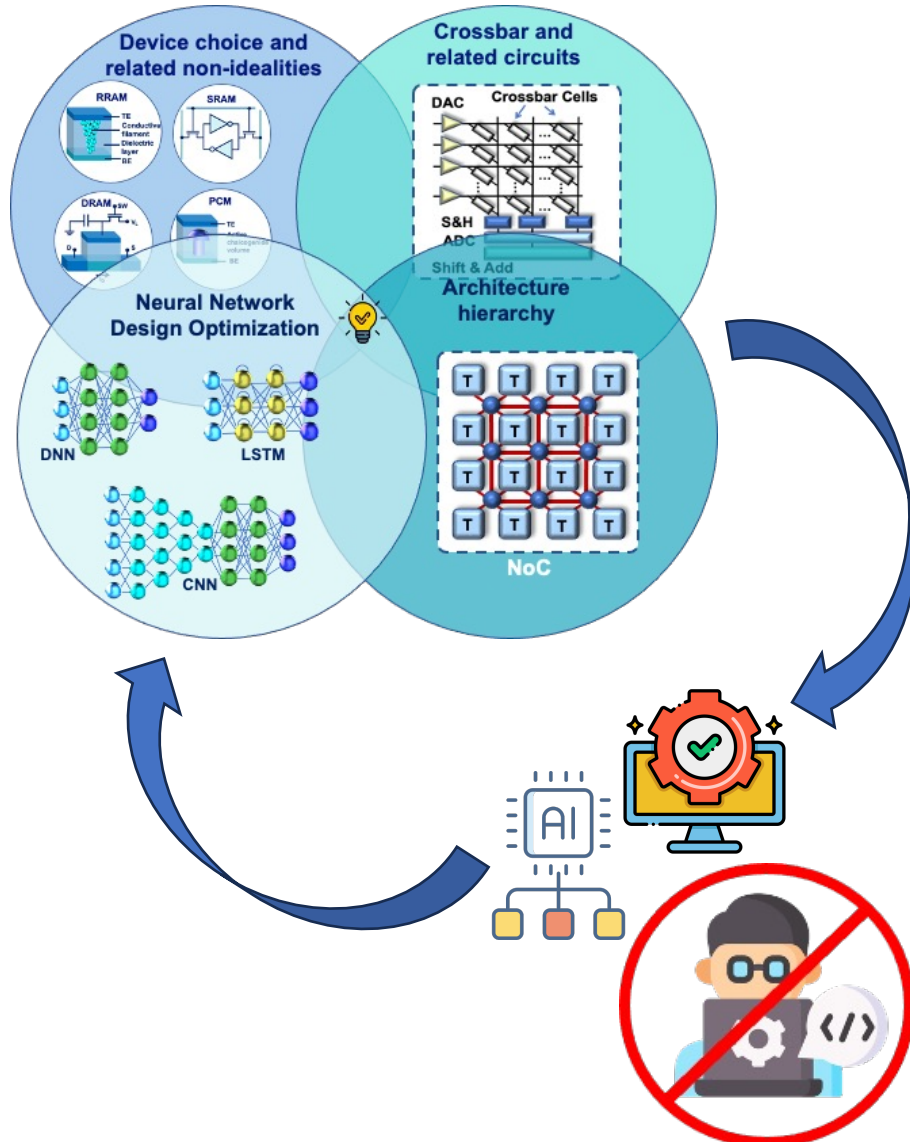
Neural network model and IMC hardware optimization covered by HW-NAS:



Combining HW-NAS with other optimization techniques for end-to-end IMC hardware optimization tool?

# One more step forward

## Optimum AI hardware design



## Self-adapting Design Algorithms

## AI-driven Design Tools

- Fully-automated NAS methods capable of constructing new deep learning operations and algorithms suitable for IMC with minimal human design efforts.
- Example from software: AutoML- Zero – automatically searching for the complete machine learning algorithms (model, optimization procedure, etc with minimum restriction on the form or math operations).
- Reducing human intervention in the design.
- No pre-defined blocks.
- Adaptable to different tasks and constraints.
- IMC awareness – open challenge.



**Utilization of AI capabilities to improve and automate both algorithm and hardware design.**



# Thank you!



[ok@ieee.org](mailto:ok@ieee.org)