

# A SYSTEMATIC AND RAPID APPROACH TO DESIGN SPACE EXPLORATION FOR TENSOR ACCELERATORS

QIJING JENNY HUANG, NVIDIA  
[jennyhuang@nvidia.com](mailto:jennyhuang@nvidia.com)

\* The views expressed in this presentation are those of the speakers and do not necessarily reflect the views or positions of any entities they represent.



# Outline

1. Accelerator design space exploration (DSE)
  2. Taxonomy of DSE Tools
  3. An overview of our approach
    - An optimization-driven mapper: CoSA
    - A search space transformation: VAESA
    - A differentiable formulation: DOSA
  4. Challenges and opportunities
- 

Performance feedback





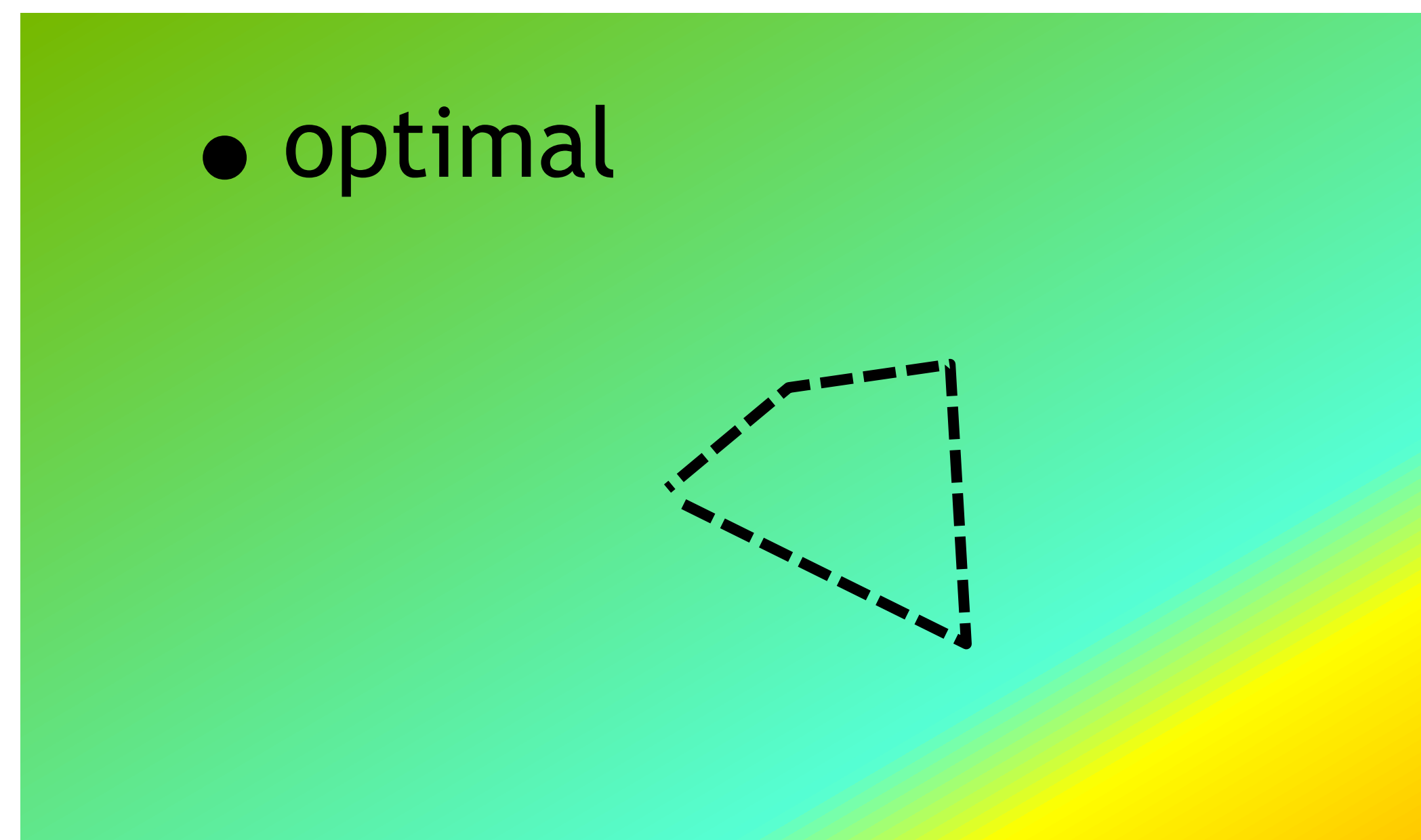
“Design is not just what it looks like  
and feels like. Design is how it  
works.”

Steve Jobs



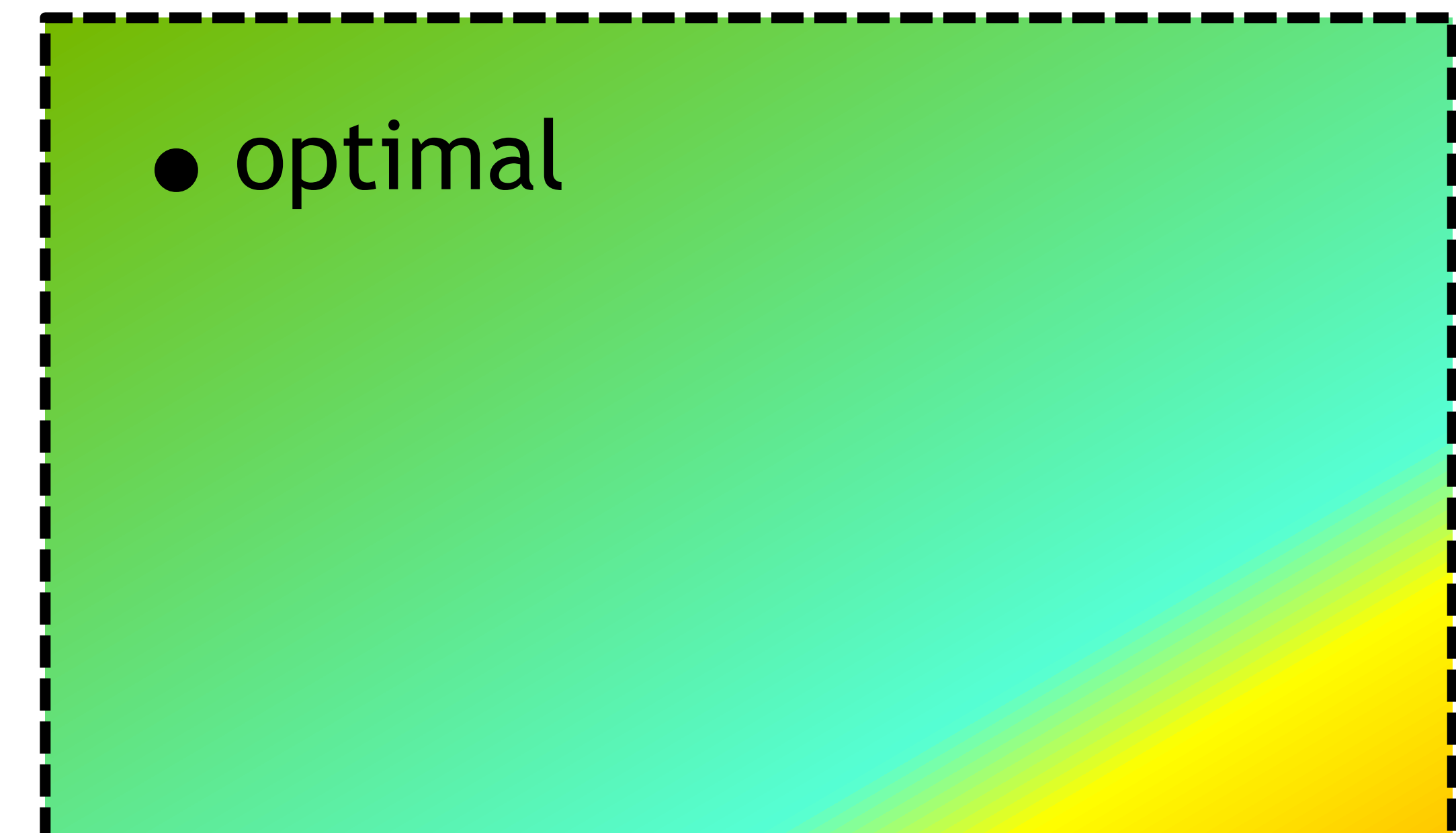
# DSE-DRIVEN ARCHITECTURE DESIGN

## Heuristics



Design Space

## DSE



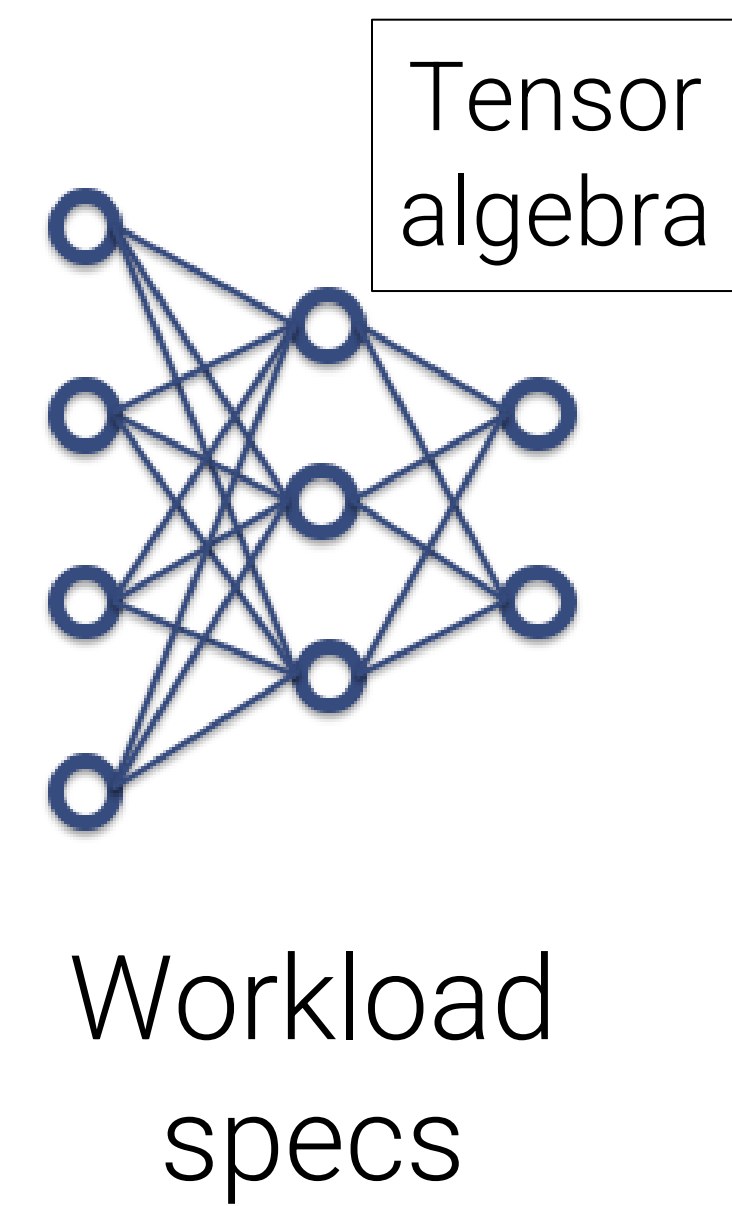
Design Space

- Architectural design should be done using DSE with:
  - Clearly defined objectives
  - A vast design space

# ACCELERATOR DESIGN SPACE EXPLORATION

Four key steps

**Step #1: Define the design space and the objectives**



# TARGET WORKLOADS

## Tensor Algebra

- **Tensor algebra** is a category of computation that can be expressed by symbols and operations of tensors
  - Example workloads:
    - Matrix-Matrix Mult, Conv, BLAS, ...

- Algebraic expression:

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

- Implementation:

```
for i in [0, I):  
  for j in [0, J):  
    for k in [0, K):  
      C[i][j] +=  
        A[i][k] * B[k][j]
```

Matrix-Matrix Mult

# TARGET WORKLOADS

## Tensor Algebra

- **Tensor Algebra** is a category of computation that can be expressed by symbols and operations of tensors
  - Example workloads:
    - Matrix-Matrix Mult, Conv, BLAS, ...

- Algebraic expression:

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

- Implementation:

```
for i in [0, I):  
  for j in [0, J):  
    for k in [0, K):  
      C[i][j] +=  
        A[i][k] * B[k][j]
```

Matrix-Matrix Mult

### Properties

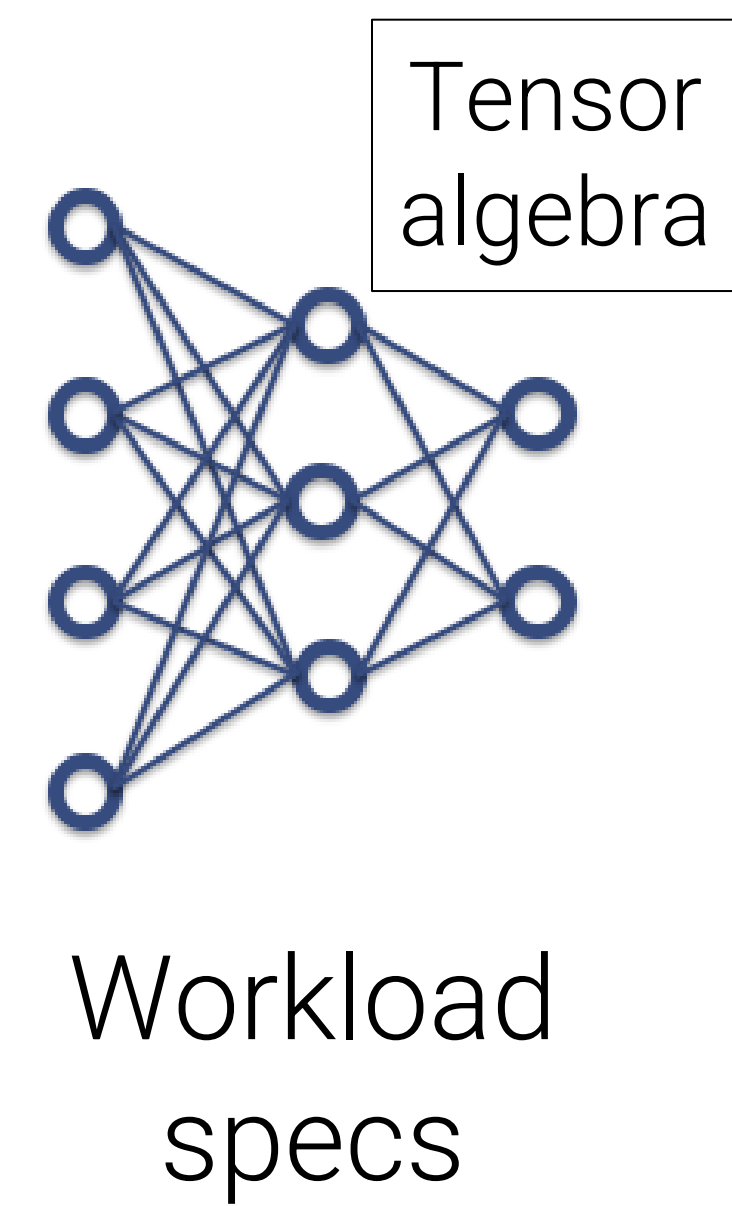
1. Known iteration space bounds
2. Regular memory access patterns
3. Repeated control flow

These properties give rise to many optimizations in accelerator DSE

# ACCELERATOR DESIGN SPACE EXPLORATION

Four key steps

**Step #1: Define the design space and the objectives**

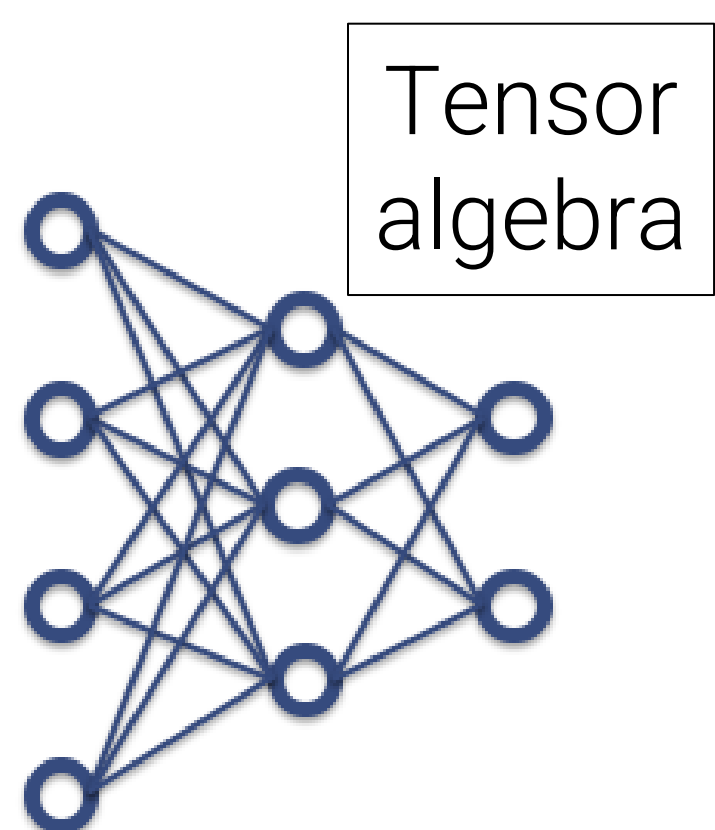




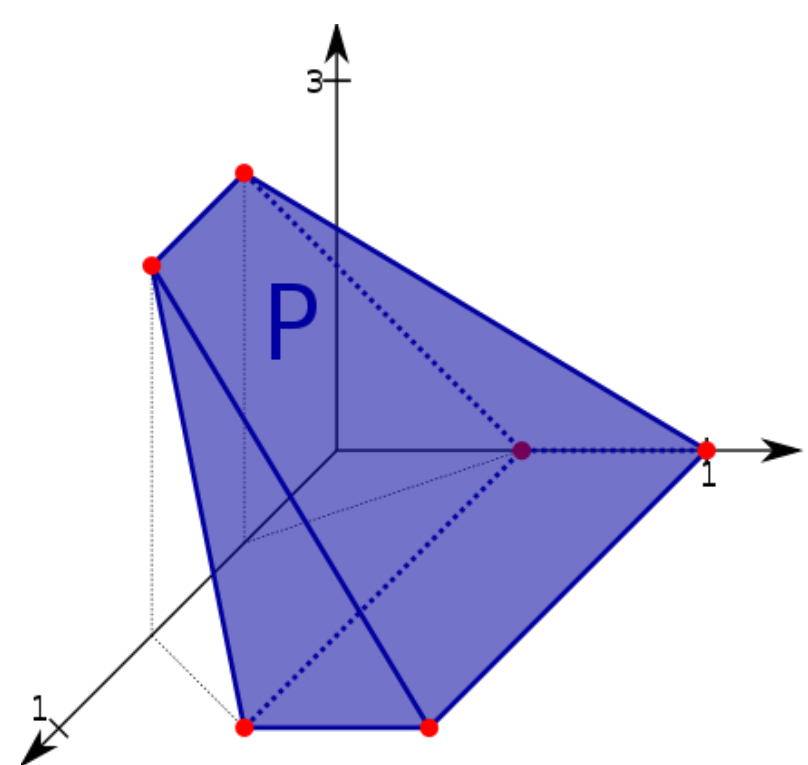
# ACCELERATOR DESIGN SPACE EXPLORATION

Four key steps

## Step #1: Define the design space and objectives



Workload  
specs



Mapping  
constraints

Metrics
Latency
Energy
Area
EDP
...

Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

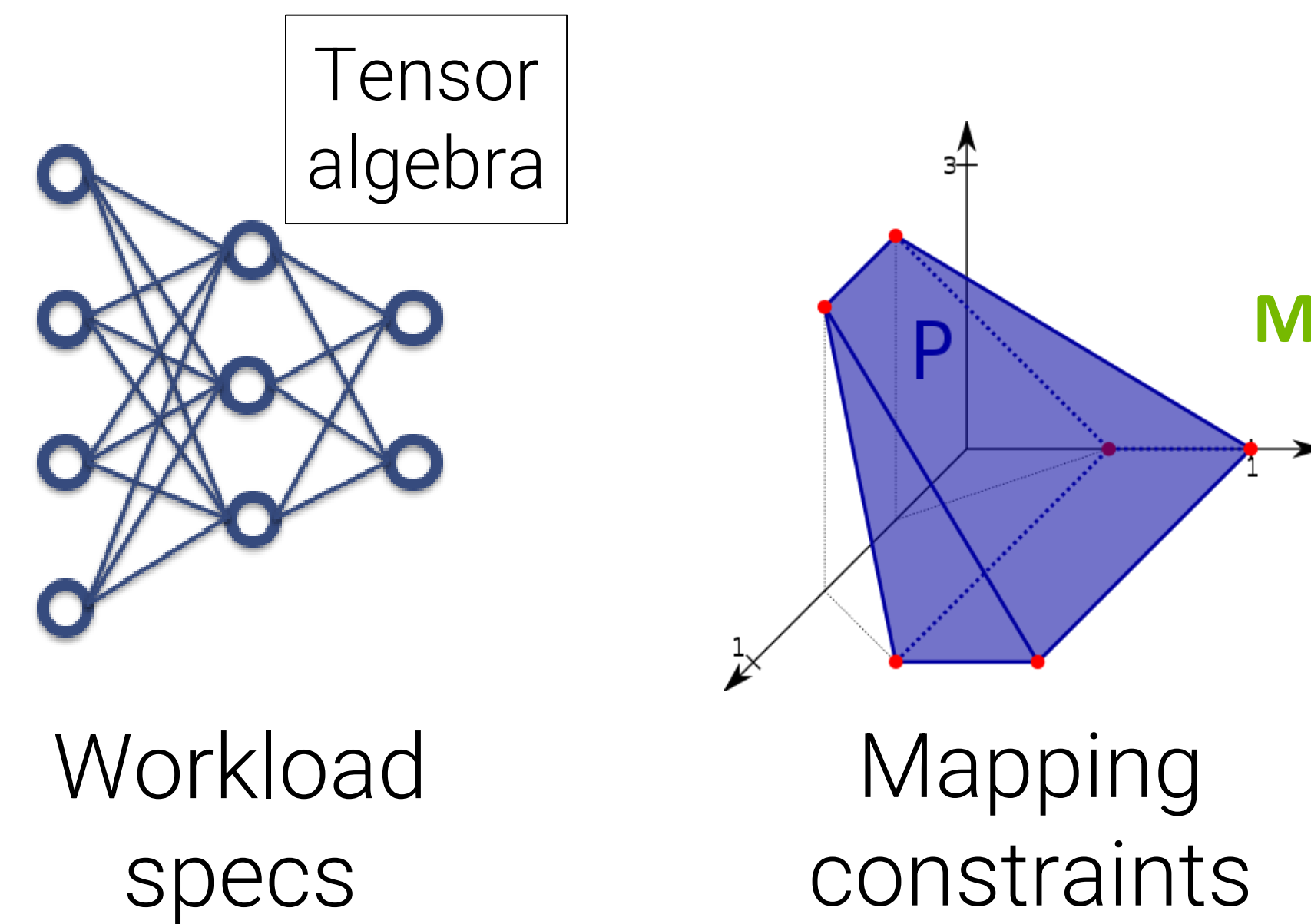
Arch design space



# ACCELERATOR DESIGN SPACE EXPLORATION

## Four key steps

### Step #1: Define the design space and objectives



Metrics
Latency
Energy
Area
EDP
...

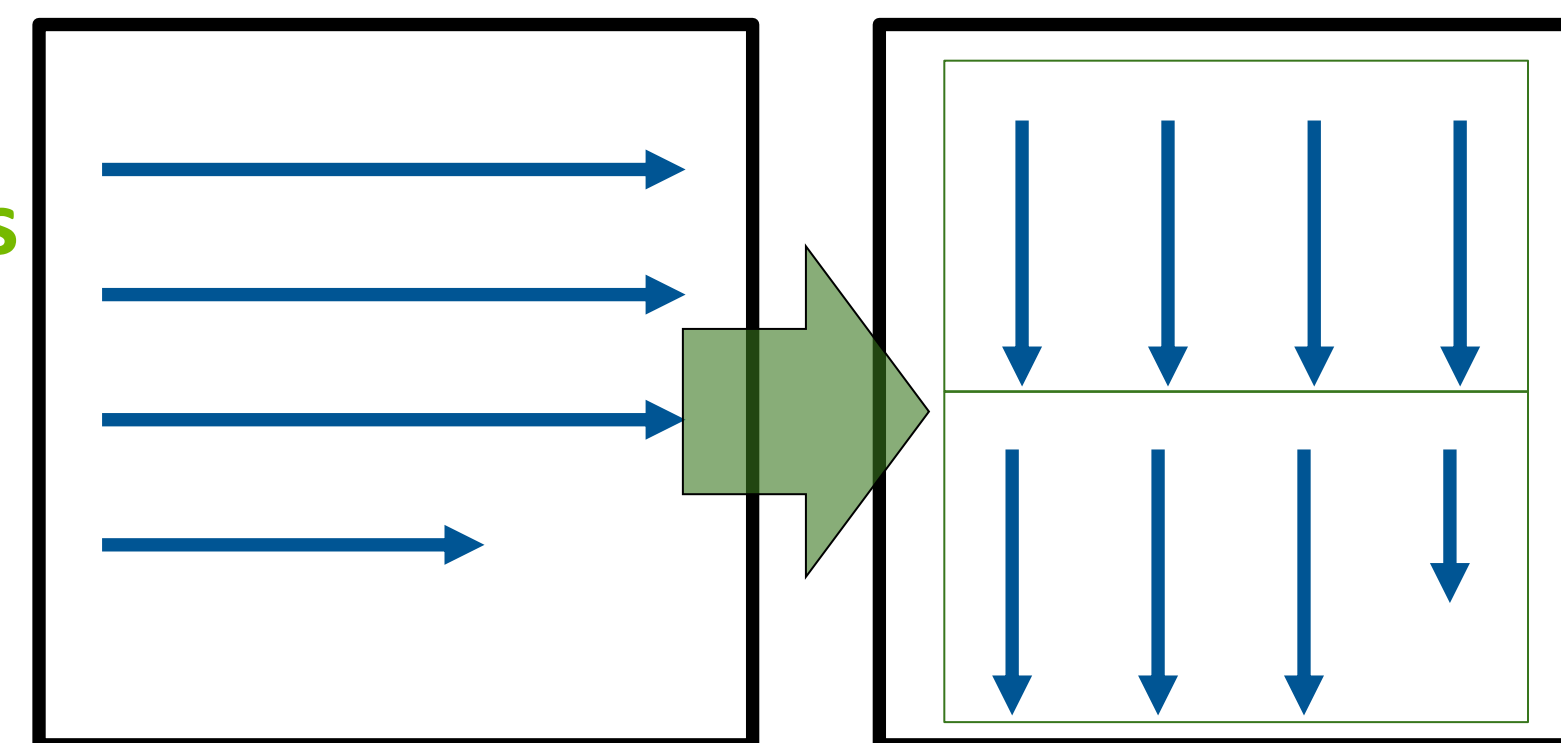
Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

Arch design space

### STEP #2: Optimize the mapping of the workload on the HW design

Mapper



- Tiling factors
- Spatial / temporal mapping
- Loop permutation

New arch configs

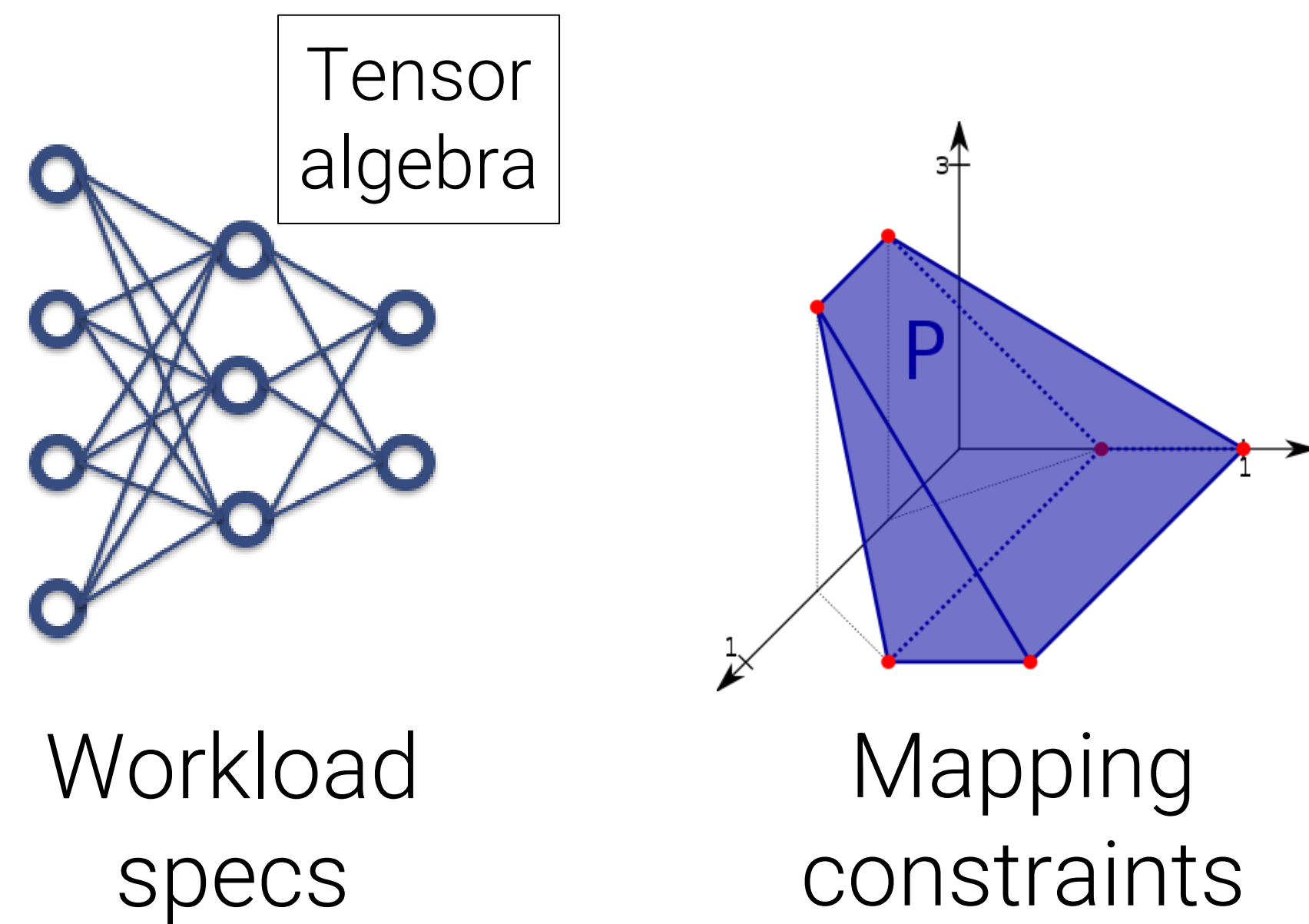
Arch design space



# ACCELERATOR DESIGN SPACE EXPLORATION

## Four key steps

### Step #1: Define the design space and objectives



Metrics
Latency
Energy
Area
EDP
...

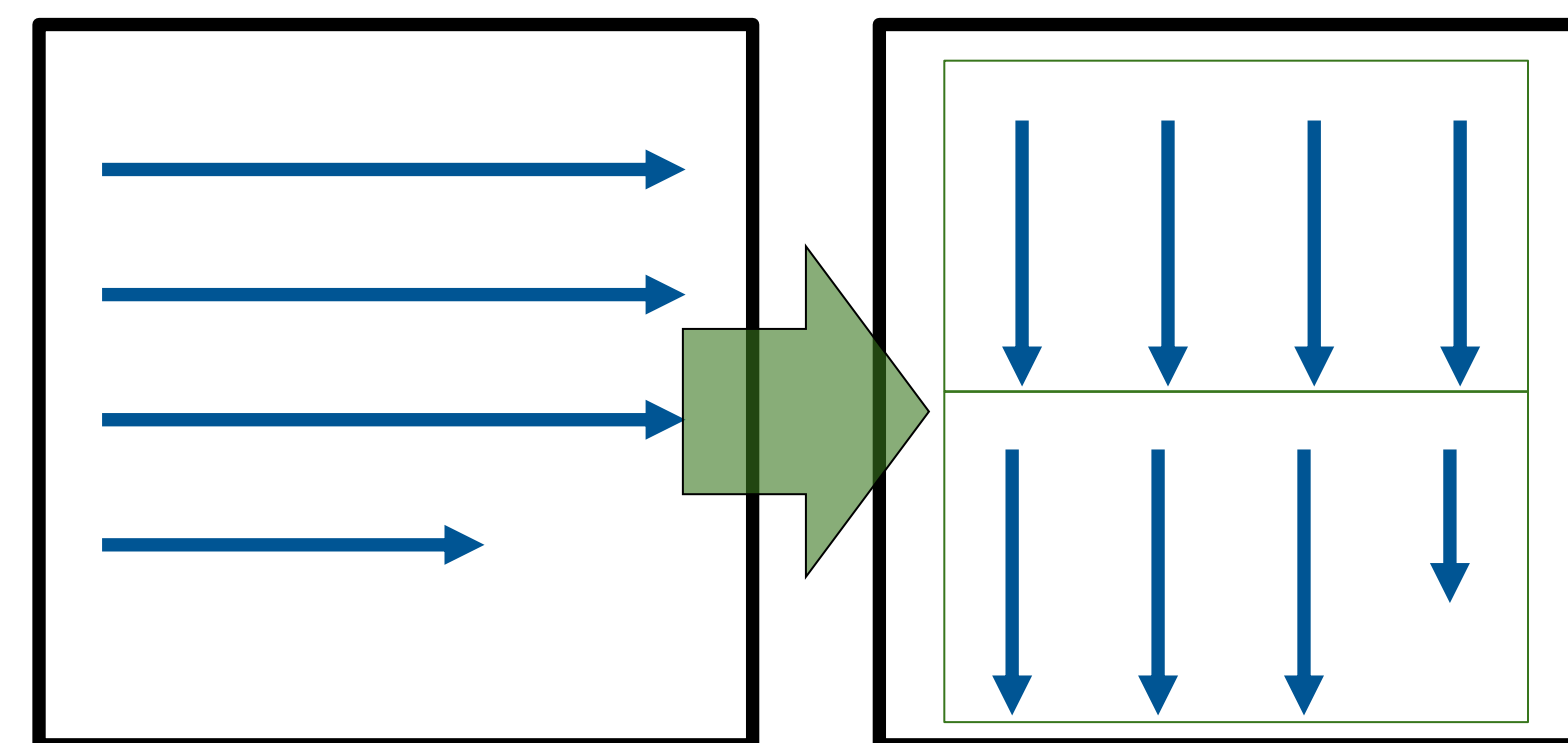
Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

Arch design space

### STEP #2: Optimize the mapping of the workload on the HW design

Mapper



- Tiling factors
- Spatial / temporal mapping
- Loop permutation

Workload specs  
Arch configs  
Objectives  
Mapping

### STEP #3: Evaluate the performance

Evaluation Tool



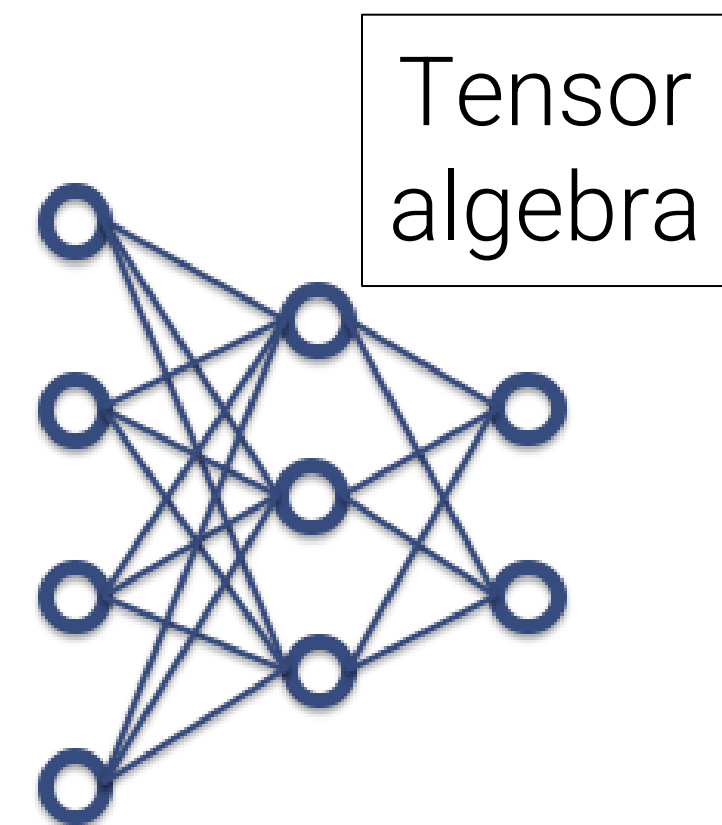
- Latency
- Energy
- Area
- ...



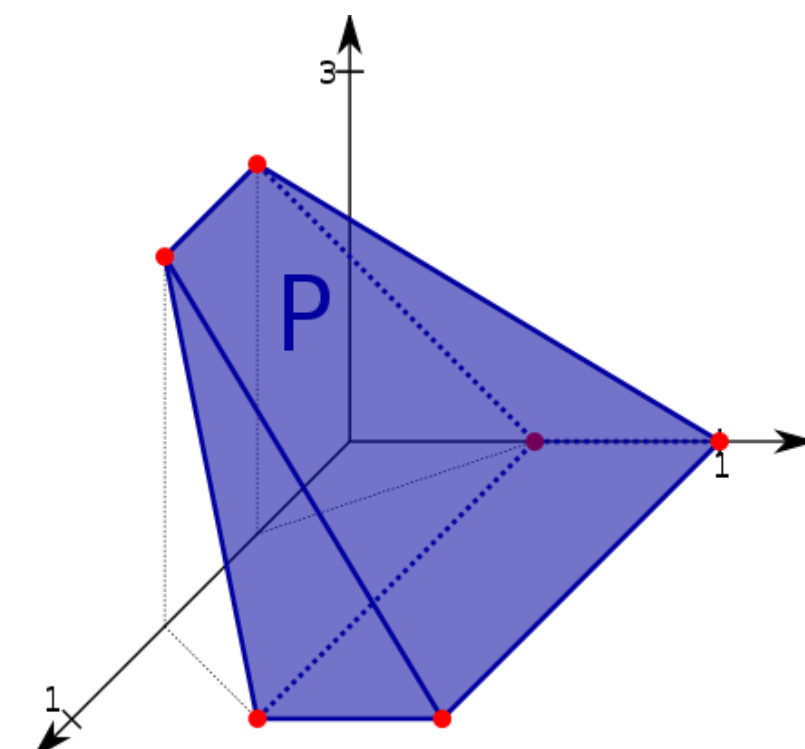
# ACCELERATOR DESIGN SPACE EXPLORATION

## Four key steps

### Step #1: Define the design space and objectives



Workload specs



Mapping constraints

Metrics
Latency
Energy
Area
EDP
...

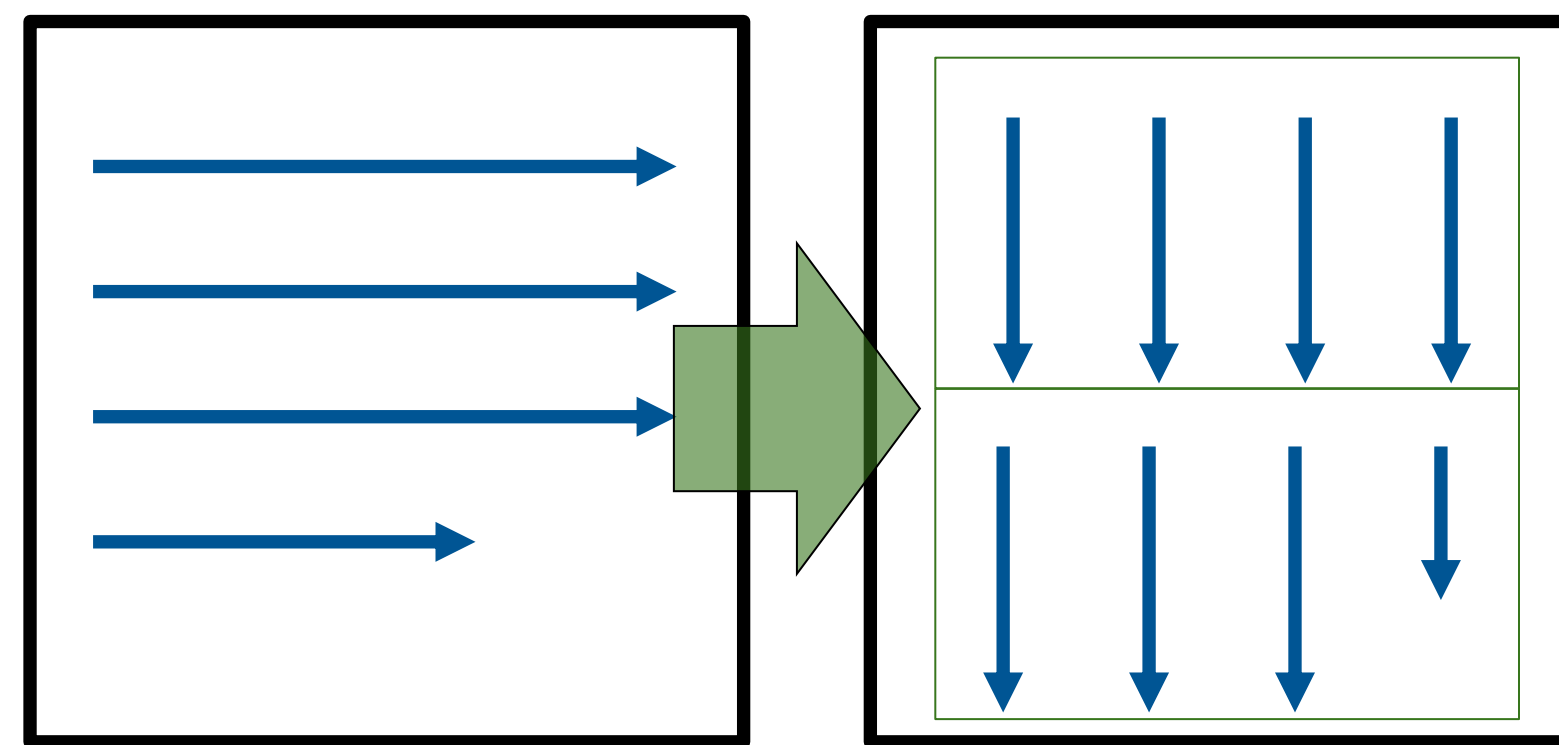
Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

Arch design space

### STEP #2: Optimize the mapping of the workload on the HW design

Mapper



- Tiling factors
- Spatial / temporal mapping
- Loop permutation

Performance feedback

### STEP #3: Evaluate the performance

Evaluation Tool



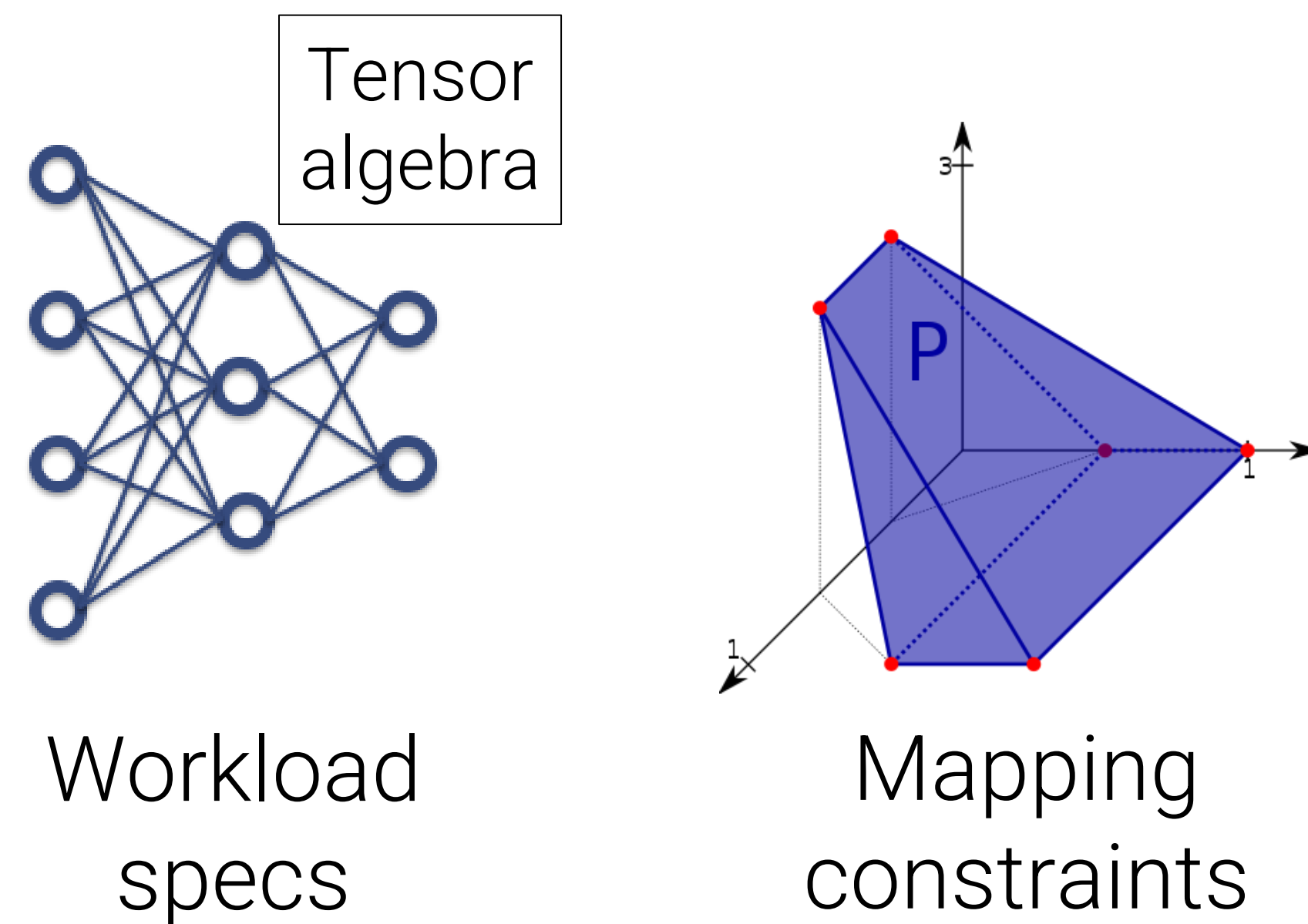
- Latency
- Energy
- Area
- ...



# ACCELERATOR DESIGN SPACE EXPLORATION

## Four key steps

### Step #1: Define the design space and objectives



Metrics
Latency
Energy
Area
EDP
...

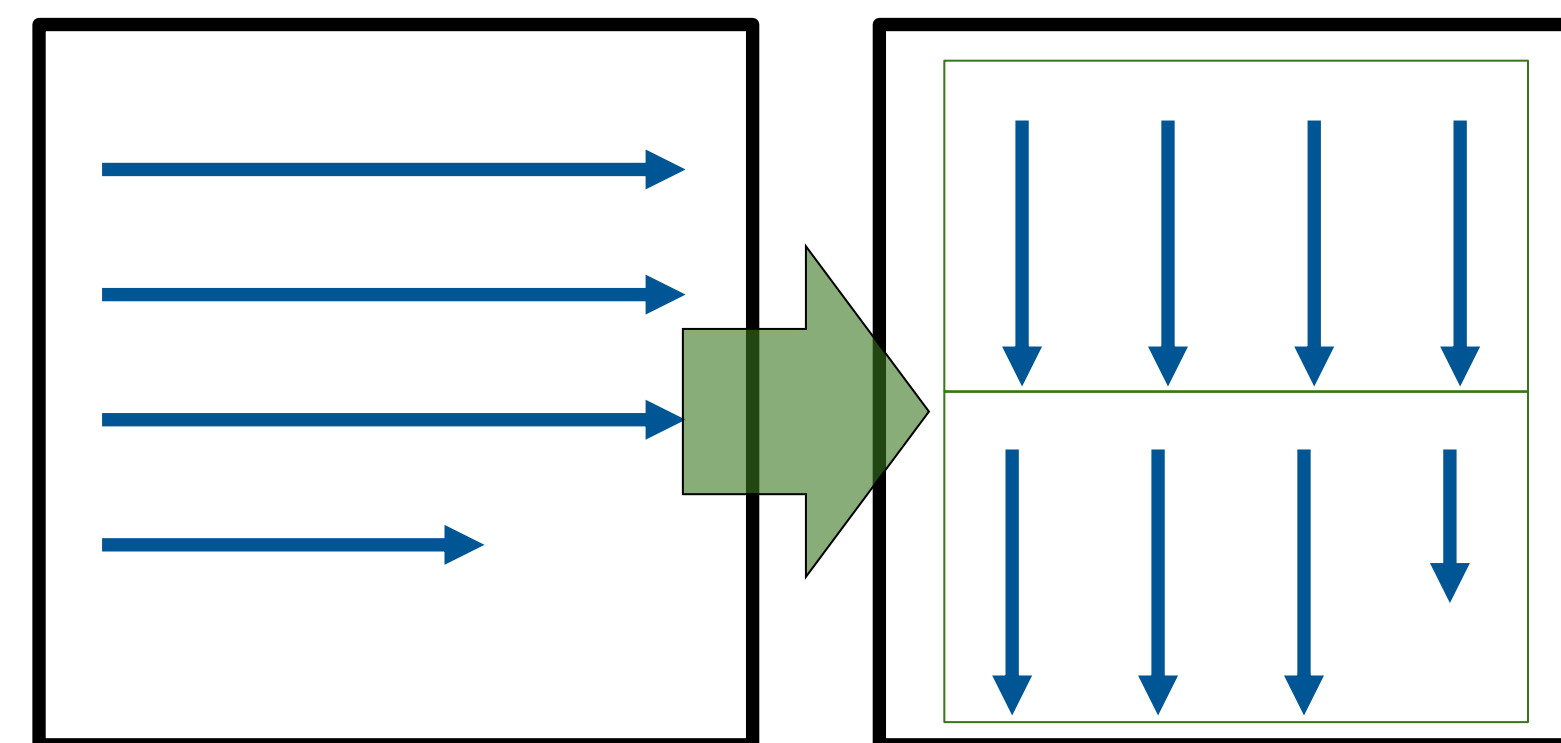
Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

Arch design space

### STEP #2: Optimize the mapping of the workload on the HW design

Mapper



- Tiling factors
- Spatial / temporal mapping
- Loop permutation

### STEP #3: Evaluate the performance

Evaluation Tool



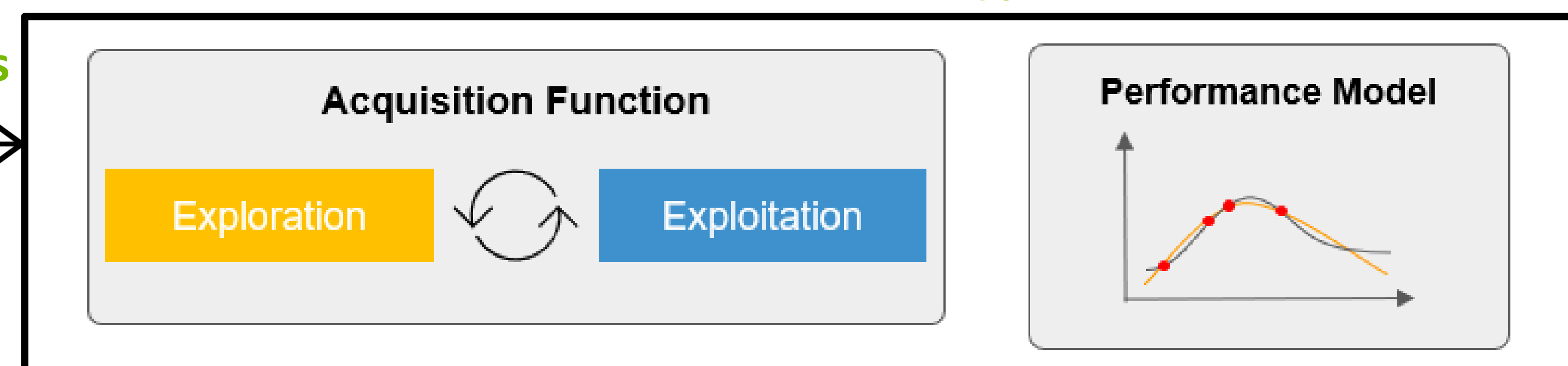
- Latency
- Energy
- Area
- ...

Performance feedback

### STEP #4: Update the search algorithm and select the next design points

Search Strategy

Workload specs  
Arch design space



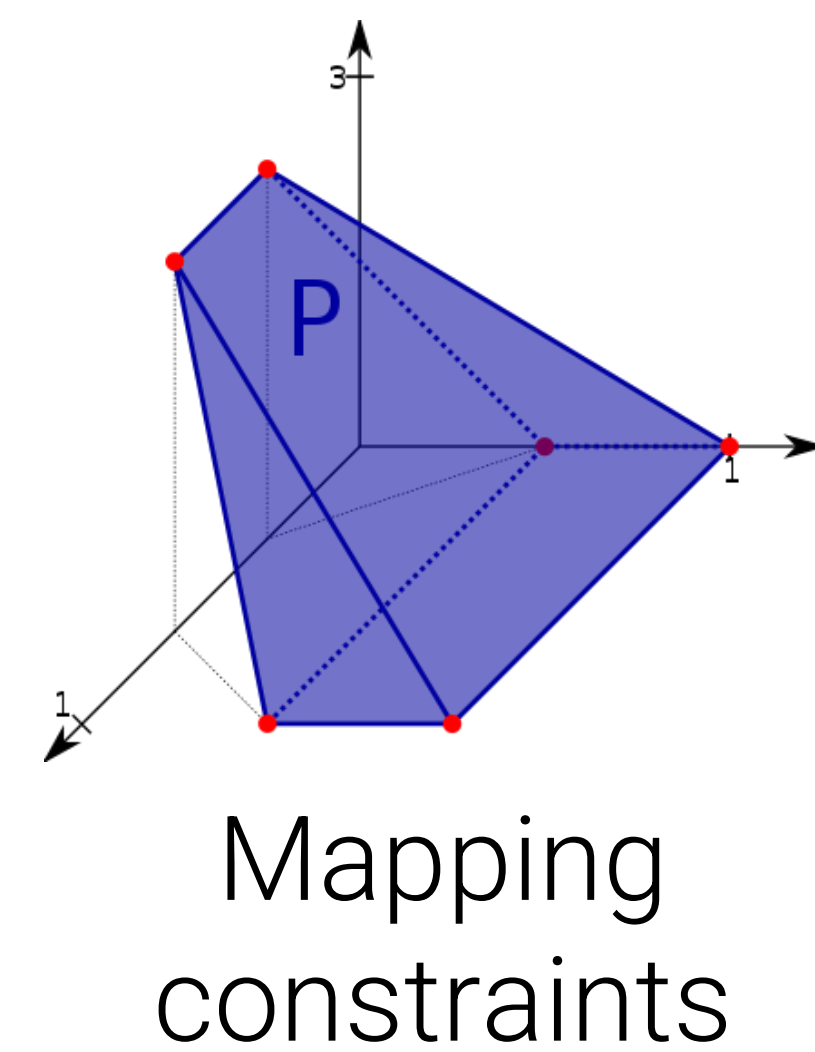
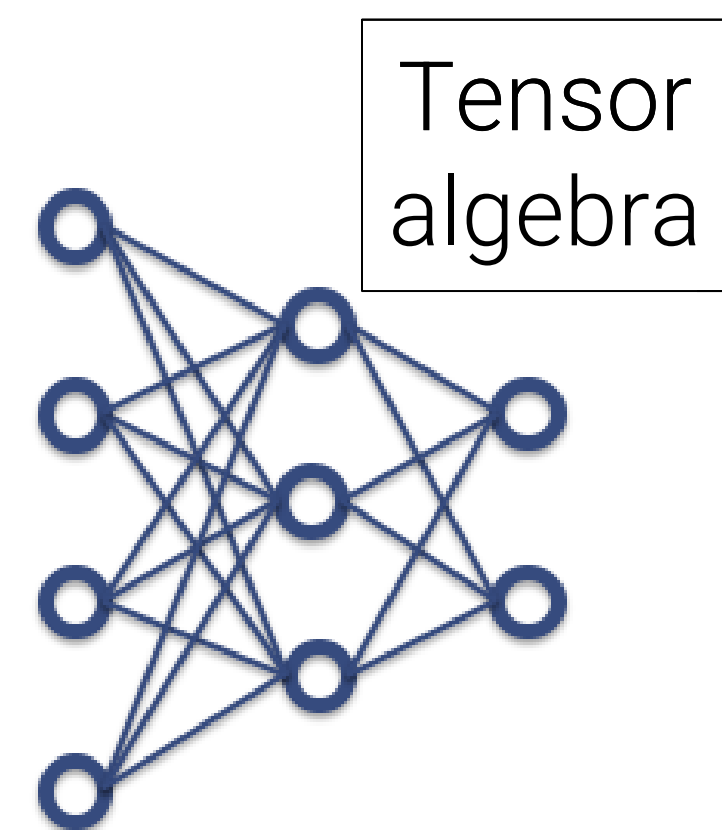
- Random Search
- Black-box Optimization
- Gradient-based Optimization
- ...



# ACCELERATOR DESIGN SPACE EXPLORATION

## Four key steps

### Step #1: Define the design space and objectives



Metrics
Latency
Energy
Area
EDP
...

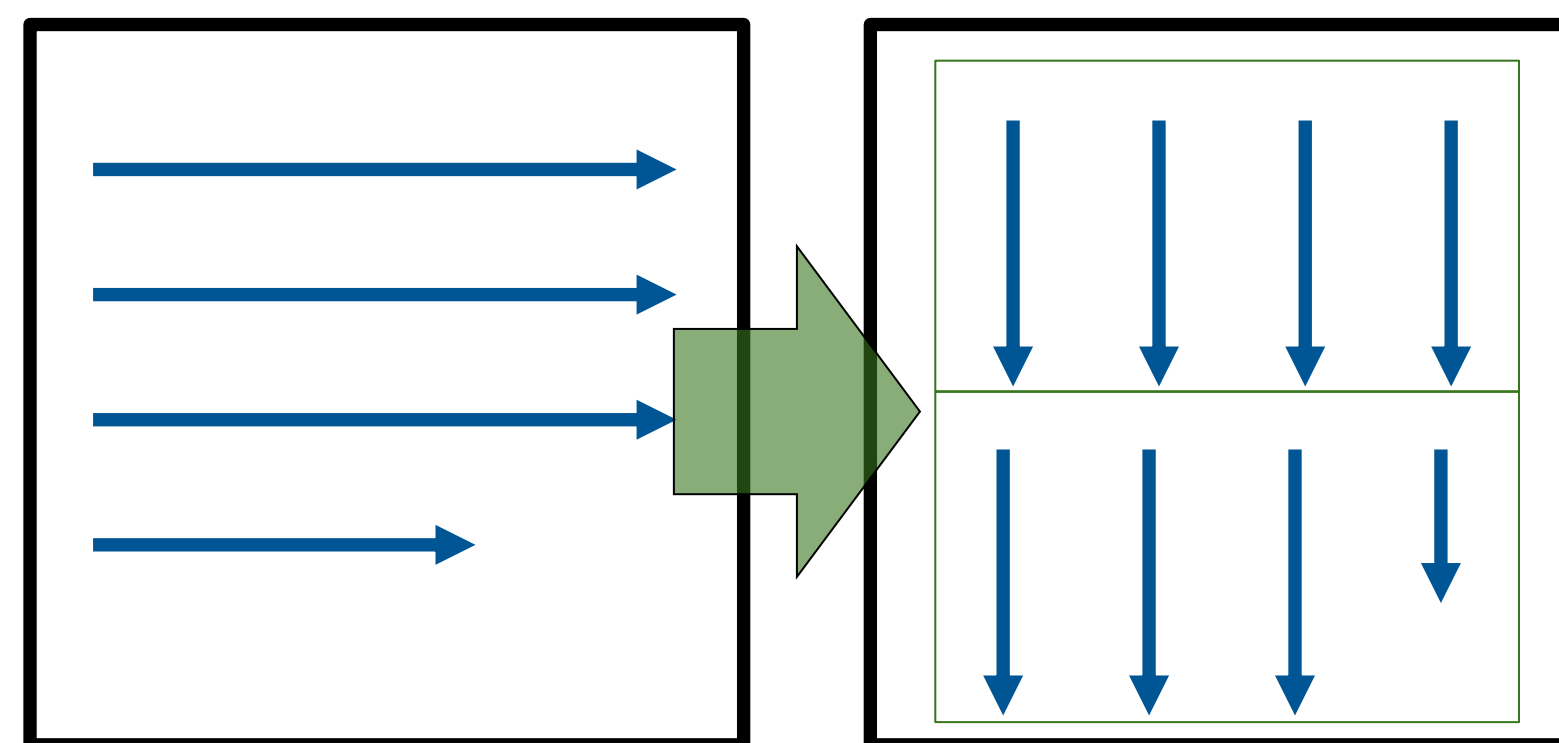
Objectives

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB

Arch design space

### STEP #2: Optimize the mapping of the workload on the HW design

Mapper



- Tiling factors
- Spatial / temporal mapping
- Loop permutation

### STEP #3: Evaluate the performance

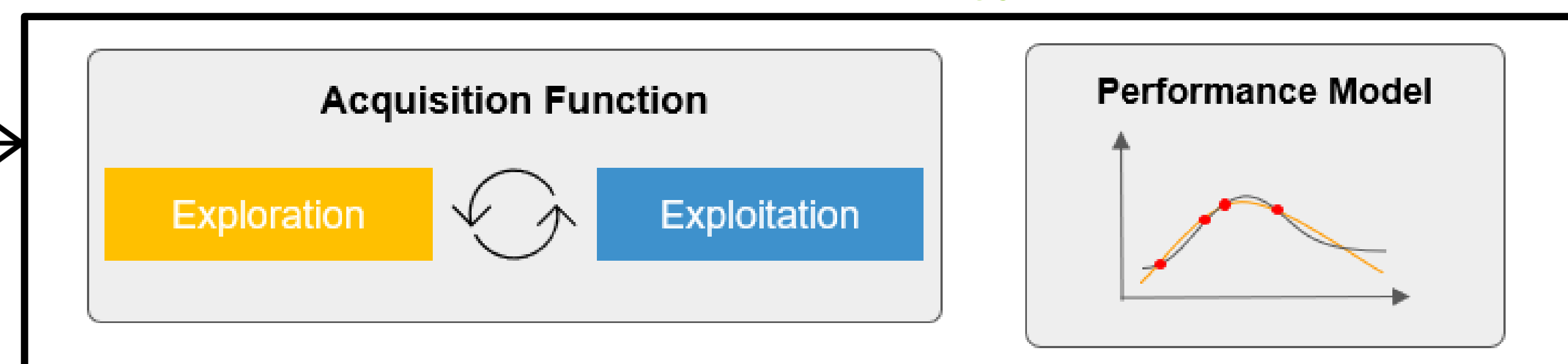
Evaluation Tool



- Latency
- Energy
- Area
- ...

### STEP #4: Update the search algorithm and select the next design points

Search Strategy



- Random Search
- Black-box Optimization
- Gradient-based Optimization
- ...



# KEY CHALLENGES IN DSE

Large design space and costly evaluation

Hardware  
Design Space

×

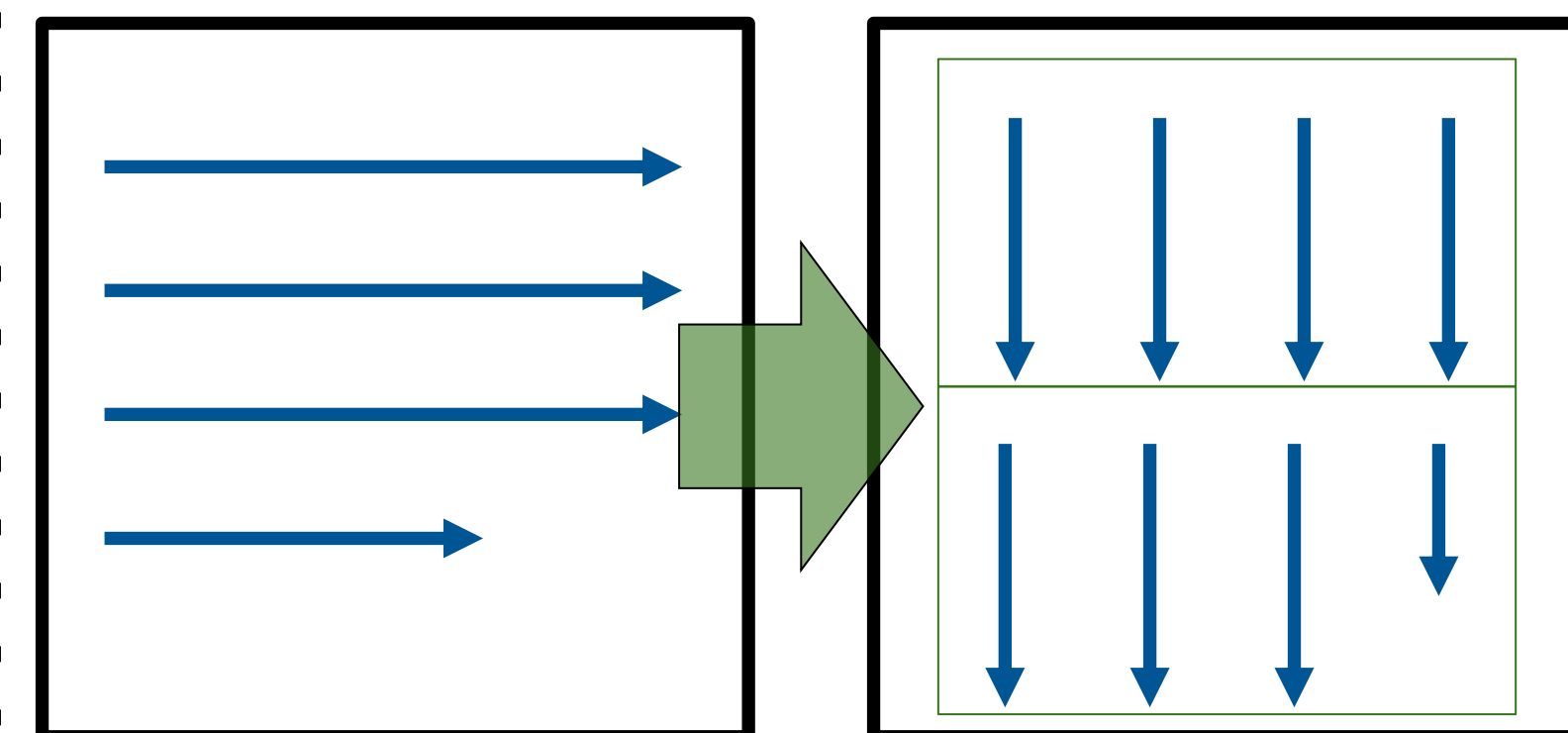
Mapping  
Space

×

Evaluation  
Time

**Intractable**

Parameter	Value Range
# of PEs	1~64
# of MAC units	1~64
Accum. buffer size	0~1MB
Weight buffer size	0~1MB
Input buffer size	0~1MB
Global buffer size	0~32MB



Platform	Evaluation Time
Timeloop	0.01s
FPGA	2 mins
VCS	10 mins
Power Analysis	6 hrs

$\sim 10^{17}$

$\sim 10^5$

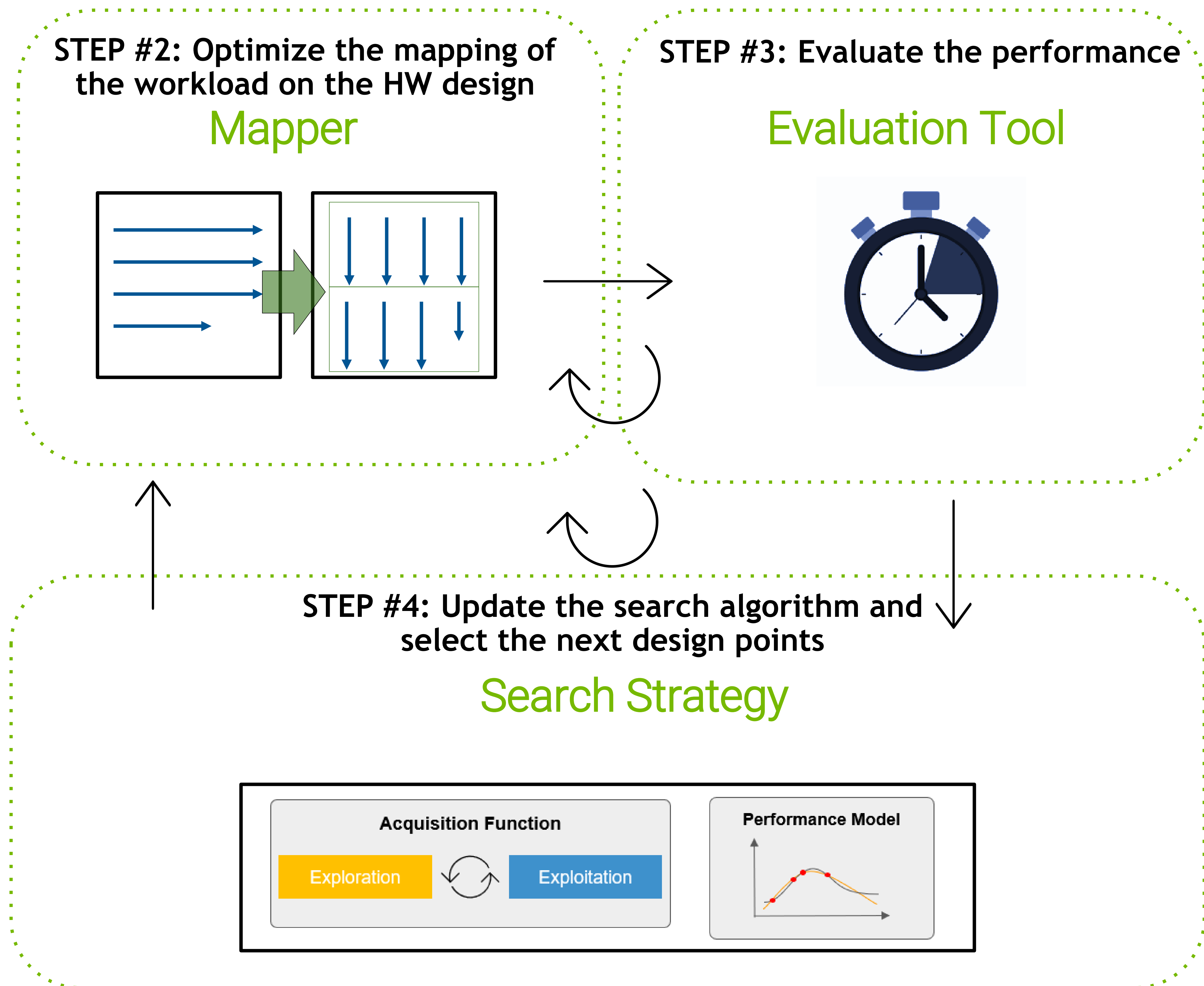
0.01s

**> 31T logical  
years**



# DESIGN SPACE EXPLORATION TOOLS

For more systematic  
and rapid DSE





# A TAXONOMY OF ACCELERATOR EVALUATION TOOLS

## Comparisons

Category	Tools	Fidelity	Modeling Speed
Polynomial Model	<a href="#">CoSA</a>	Low	Very Fast
ML Model	<a href="#">PRIME</a>	Medium	Fast†
Static Analysis	<a href="#">Timeloop</a> , <a href="#">MAESTRO</a>	Medium	Fast
Cycle-accurate Model	<a href="#">ScaleSim</a>	High	Slow
RTL Simulation	<a href="#">FireSim</a> , <a href="#">MagNet</a>	Very High	Slow*

† Varies with ML model size

\* Varies with workload size



# A TAXONOMY OF ACCELERATOR EVALUATION TOOLS

Supported features

Category	Dynamic behavior support	Data/training/implementation free	Differentiable
Polynomial Model	No	Yes	Yes
ML Model	Yes	No	Yes
Static Analysis	No	Yes	No
Cycle-accurate Model	Yes	Yes	No
RTL Simulation	Yes	No	No



# A TAXONOMY OF MAPPERS

## Heuristic-Driven

Timeloop  
Triton Marvel

- Easy to implement

## Feedback-based

AutoTVM Ansor  
Halide Gamma  
MindMapping

- More adaptive

- Costly  
- Sample invalid space  
- Hard to generalize

## Constrained Optimization

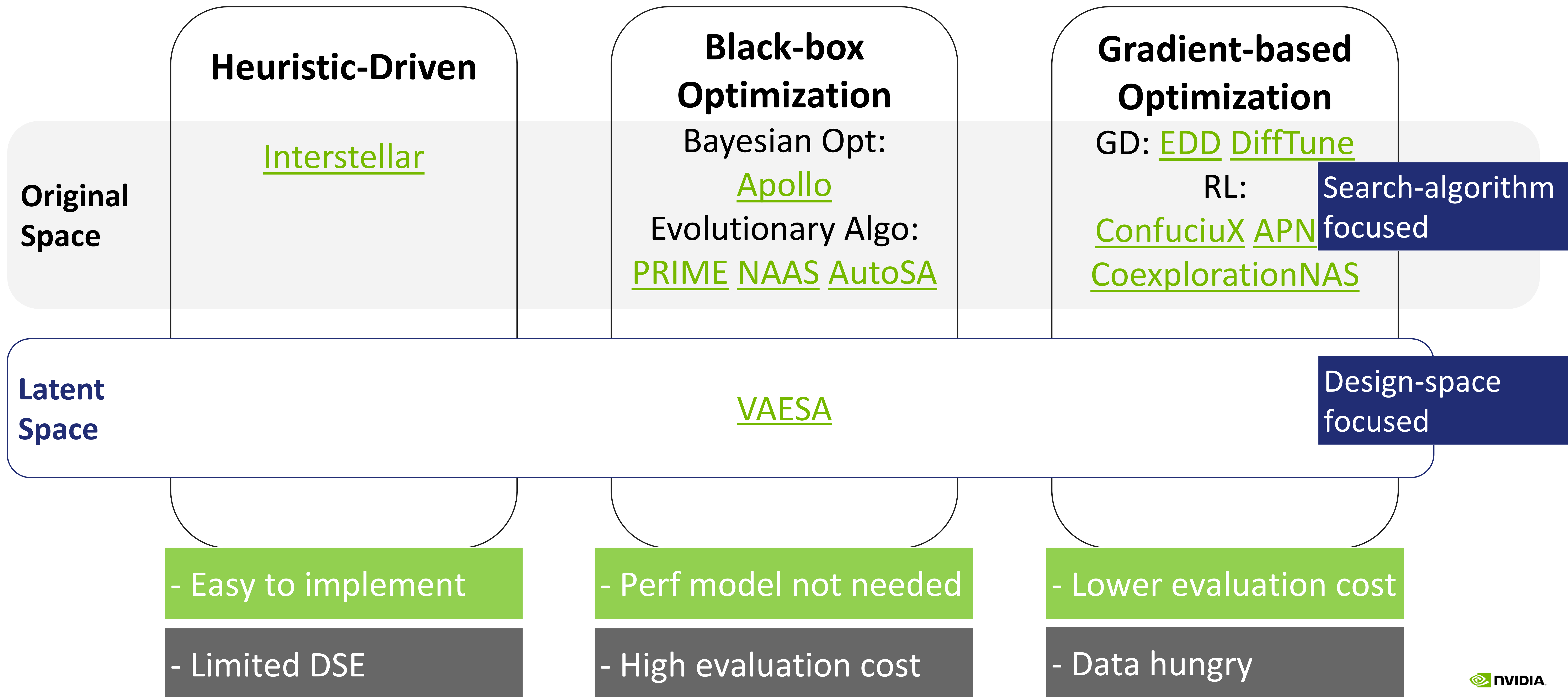
Polly+Pluto TC  
Tiramisu CoSA  
IOOpt

- More sample efficient

- Limited use case



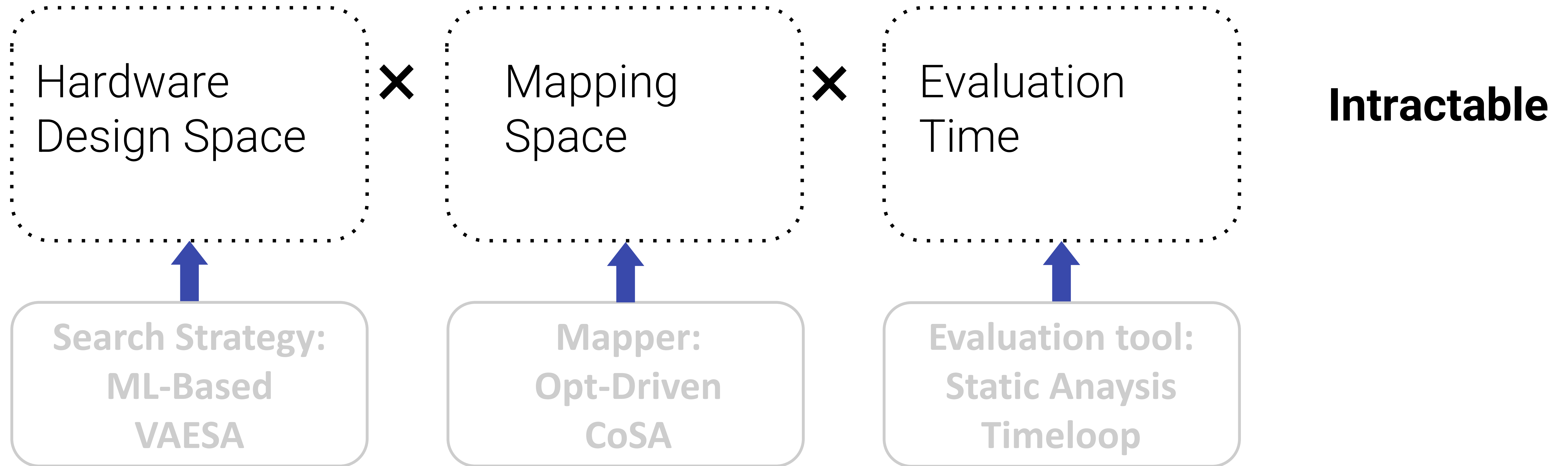
# A TAXONOMY OF ACCELERATOR SEARCH STRATEGIES





# DESIGN SPACE EXPLORATION TOOLS

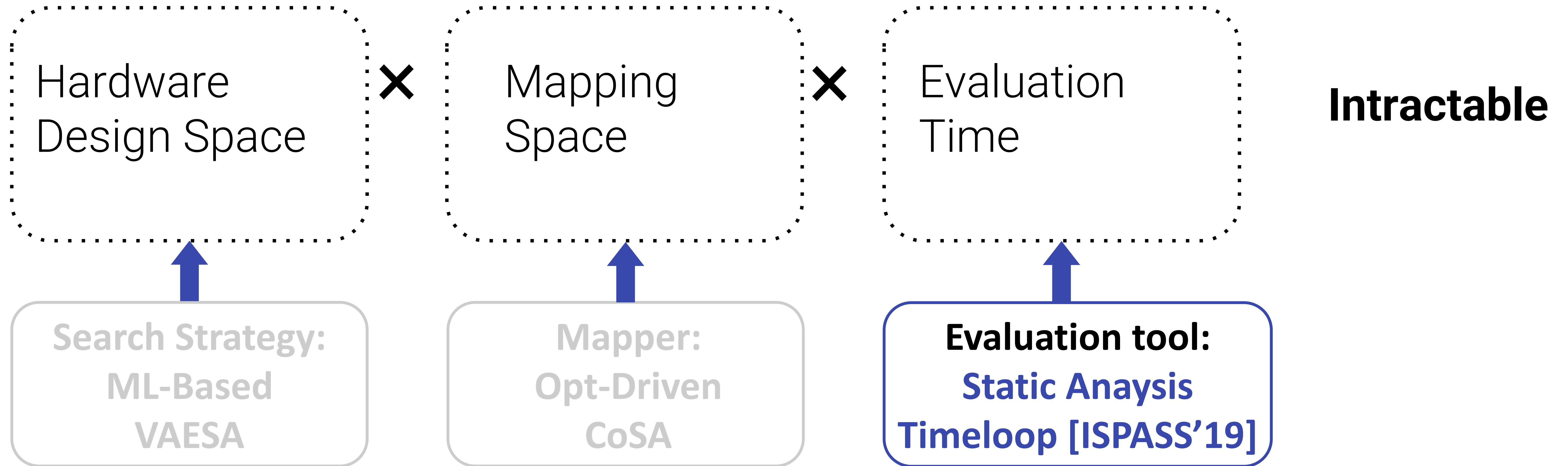
Our approach





# DESIGN SPACE EXPLORATION TOOLS

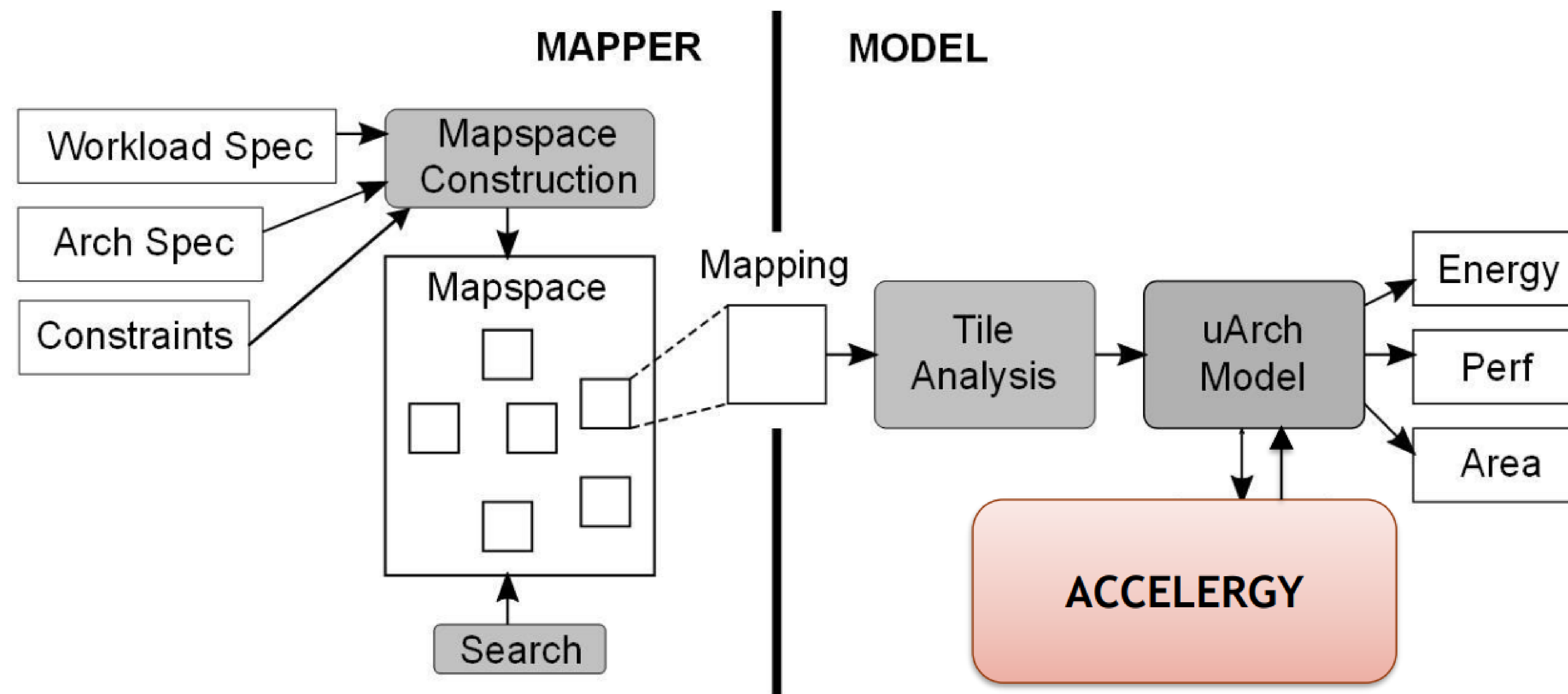
Our approach





# A STATIC ANALYSIS TOOL

Timeloop/Accelergy

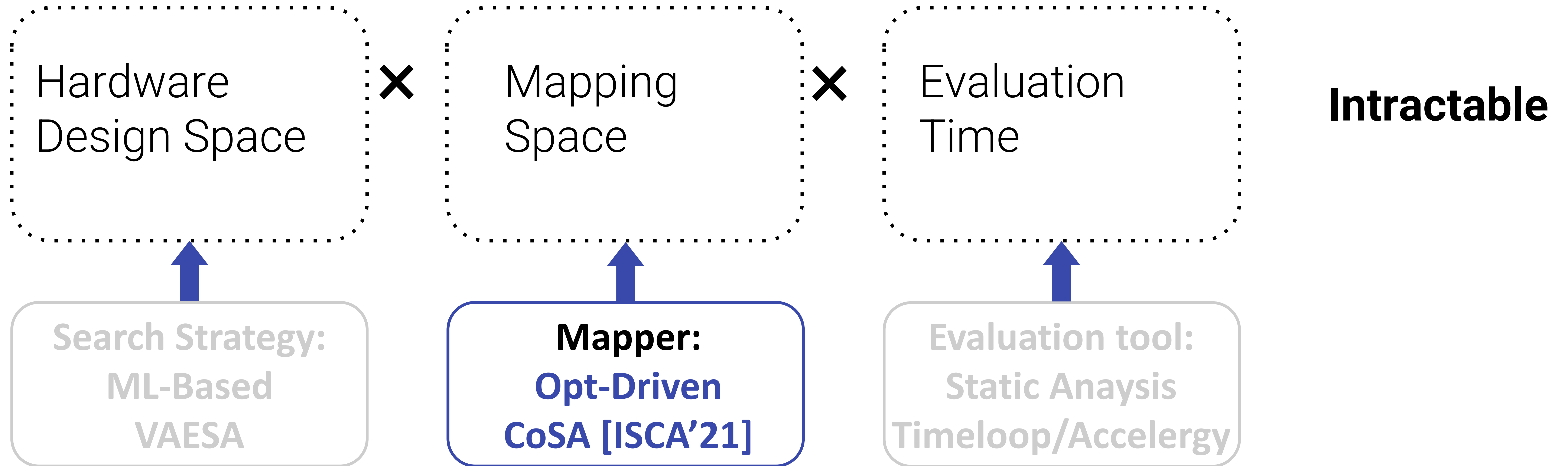


- Timeloop provides a **flexible abstraction** to define a wide range of **applications, architectures and constraints**
- Timeloop rapidly and accurately reports **latency, energy, area** using static analysis



# DESIGN SPACE EXPLORATION TOOLS

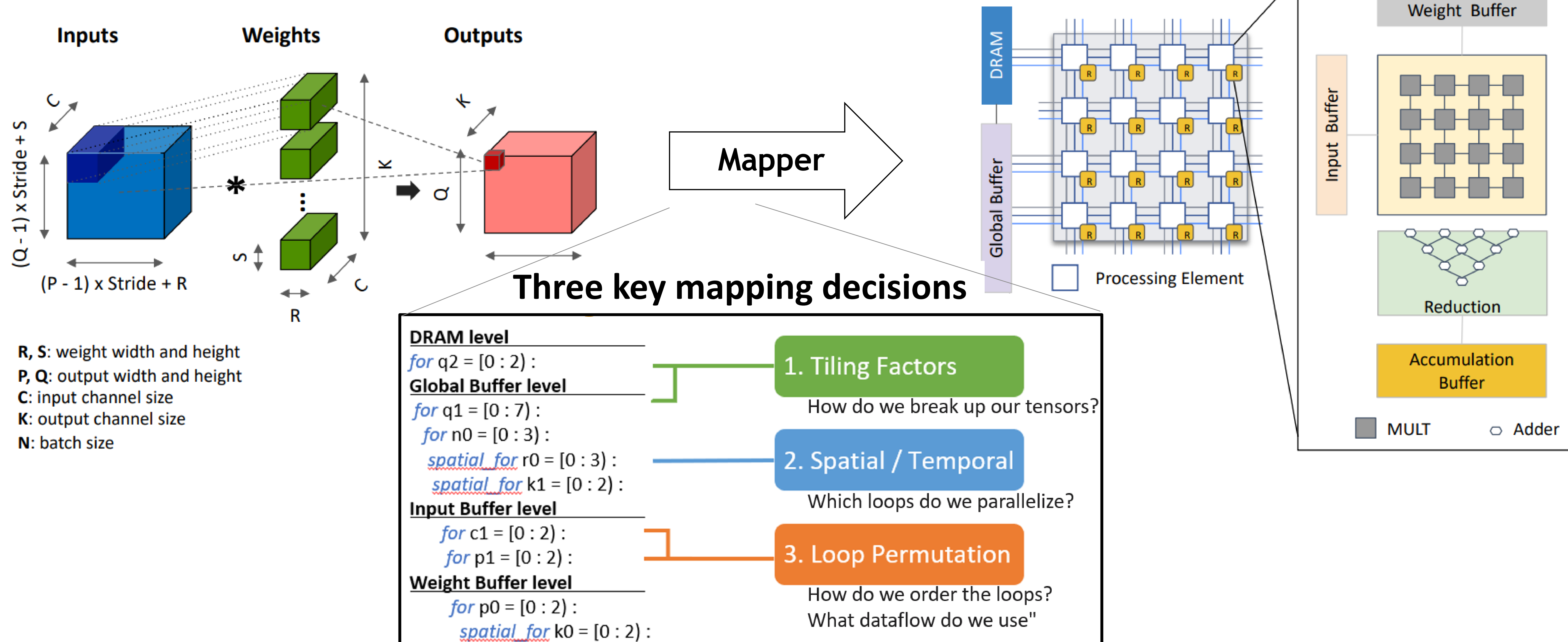
Our approach





# AN OPTIMIZATION-DRIVEN MAPPER

CoSA

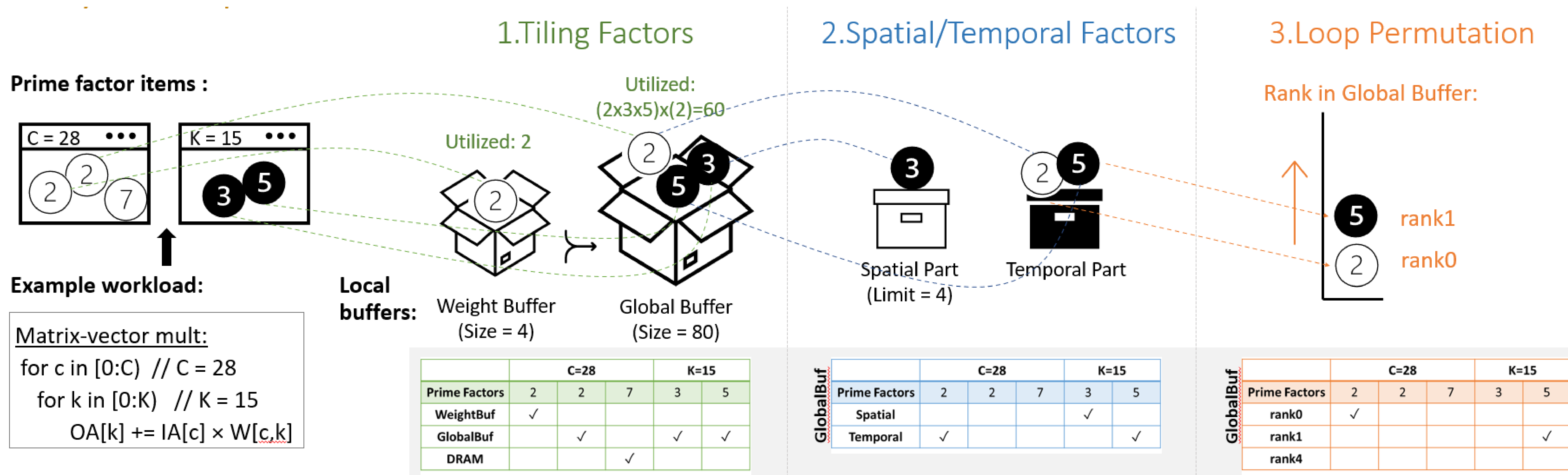


- CoSA formulates the mapping decisions into a constrained optimization problem and solves it in one shot



# AN OPTIMIZATION-DRIVEN MAPPER

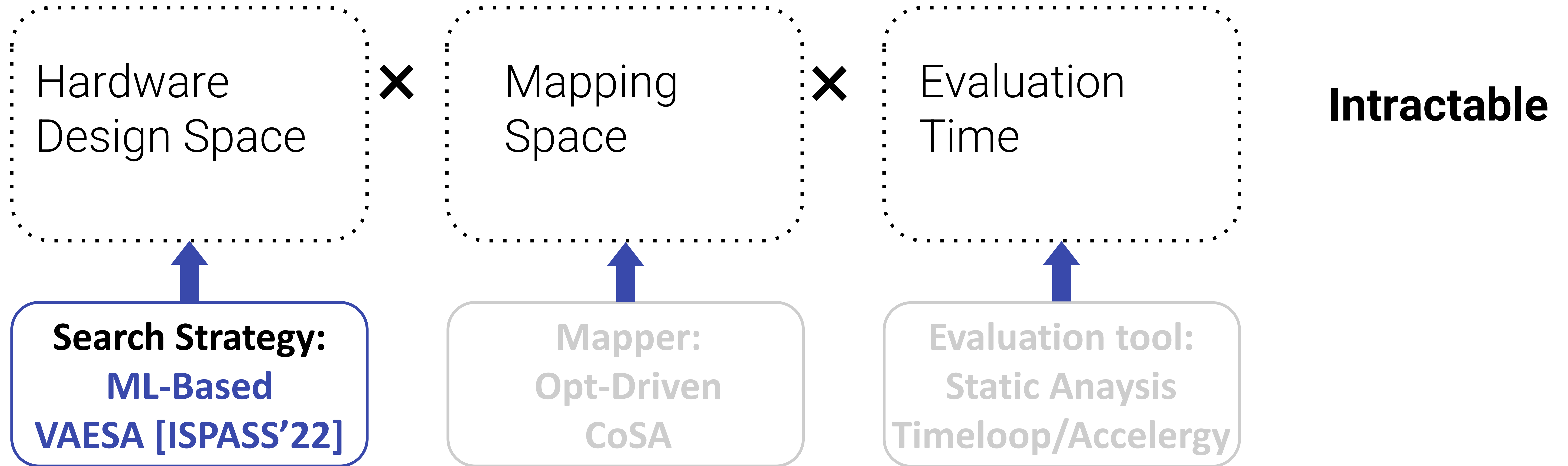
Key idea: problem factor allocation



- An optimization variable can be used to represent all three mapping decisions
- CoSA optimizes the variable using the constraints and objectives formulated in mixed integer programming
- CoSA finds mappings that are 1.5x faster and 1.2x more energy-efficient while improving the time-to-solution by 90x

# DESIGN SPACE EXPLORATION TOOLS

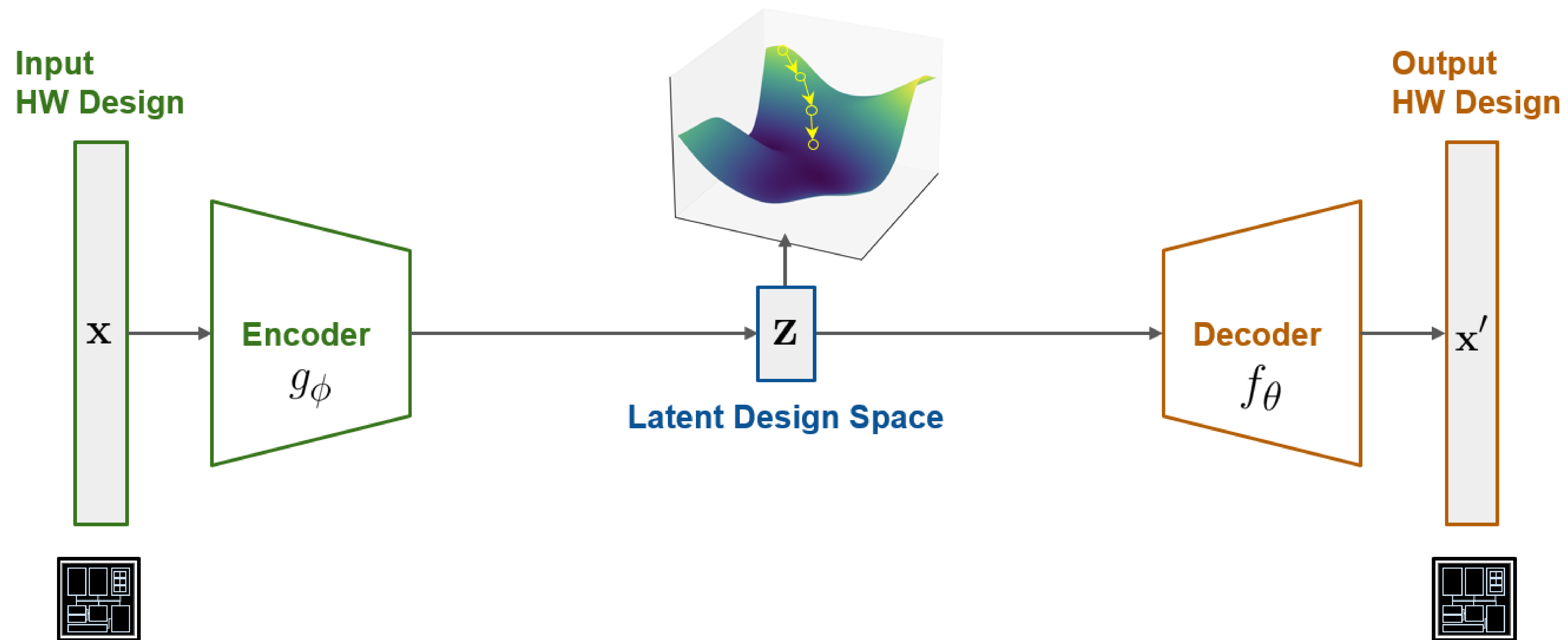
Our approach





# A ML-BASED SEARCH STRATEGY

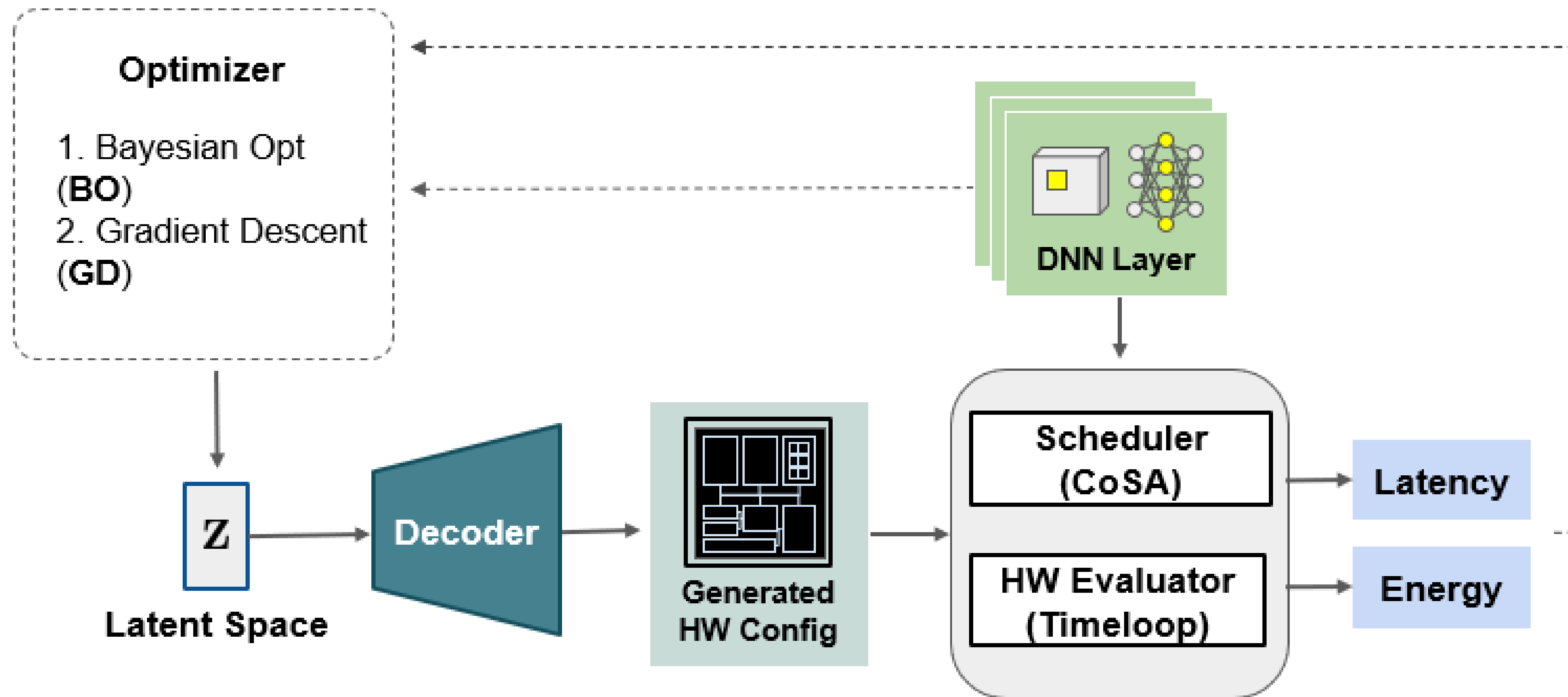
## VAESA



- **VAESA** learns a low dimensional, continuous, reconstructible latent space to facilitate accelerator DSE using Variational Autoencoder (VAE)

# A ML-BASED SEARCH STRATEGY

## VAESA Inference

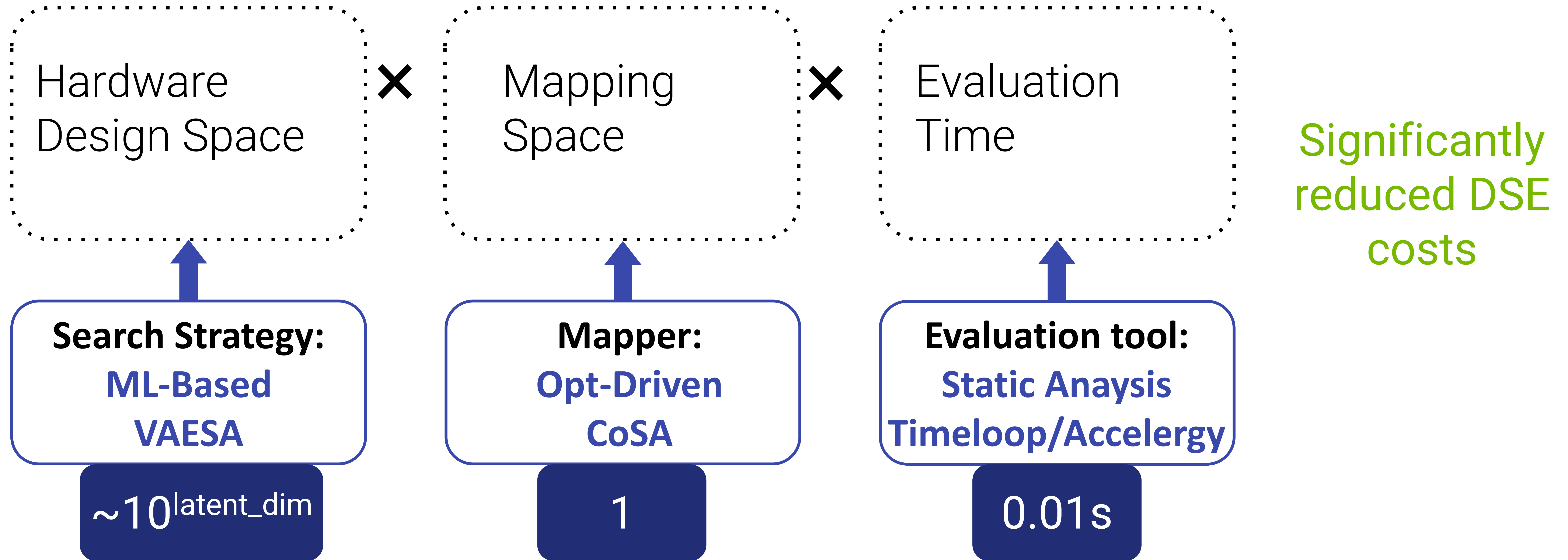


- The search algorithms are applied to the latent space and evaluated on the original search
- The latent space **reduces the search complexity** and provides a **smoother performance surface**
- Both Bayesian Optimization and Gradient Descent **achieve better sample efficiency**



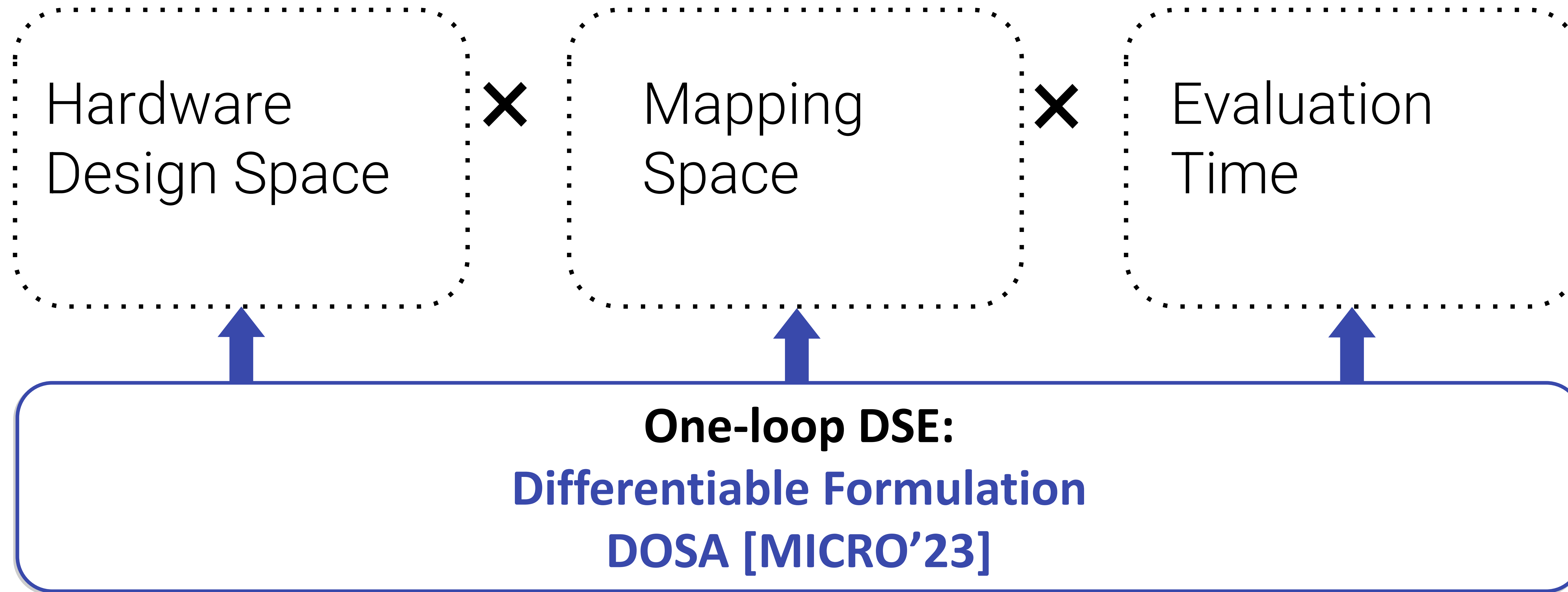
# DESIGN SPACE EXPLORATION TOOLS

Our approach



# DESIGN SPACE EXPLORATION TOOLS

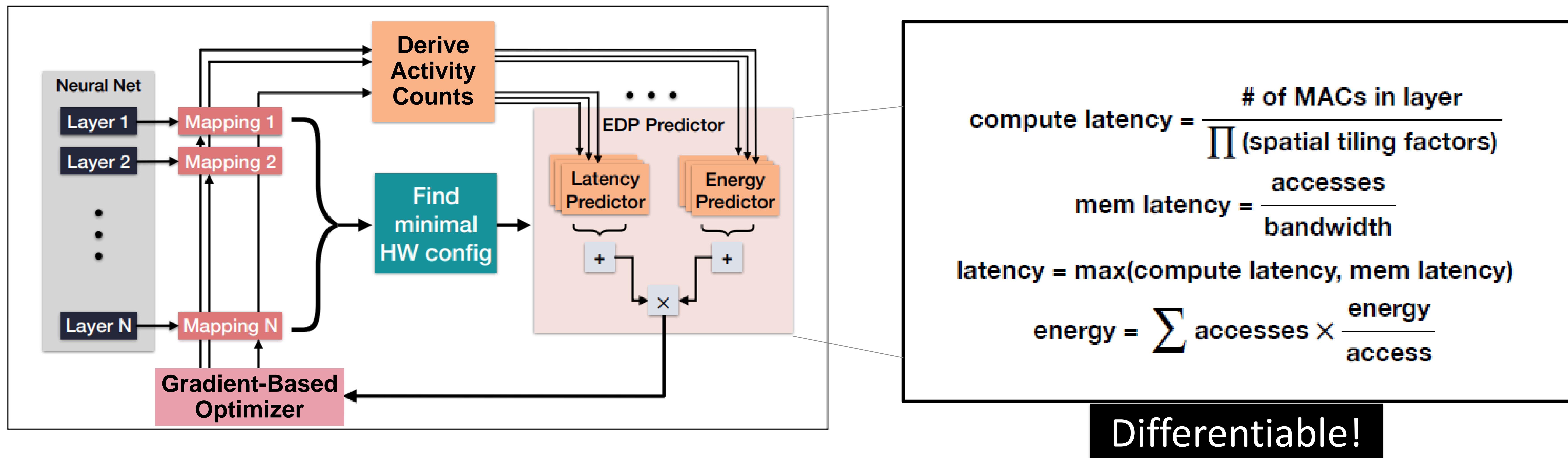
Our approach





# A DIFFERENTIABLE DSE FORMULATION

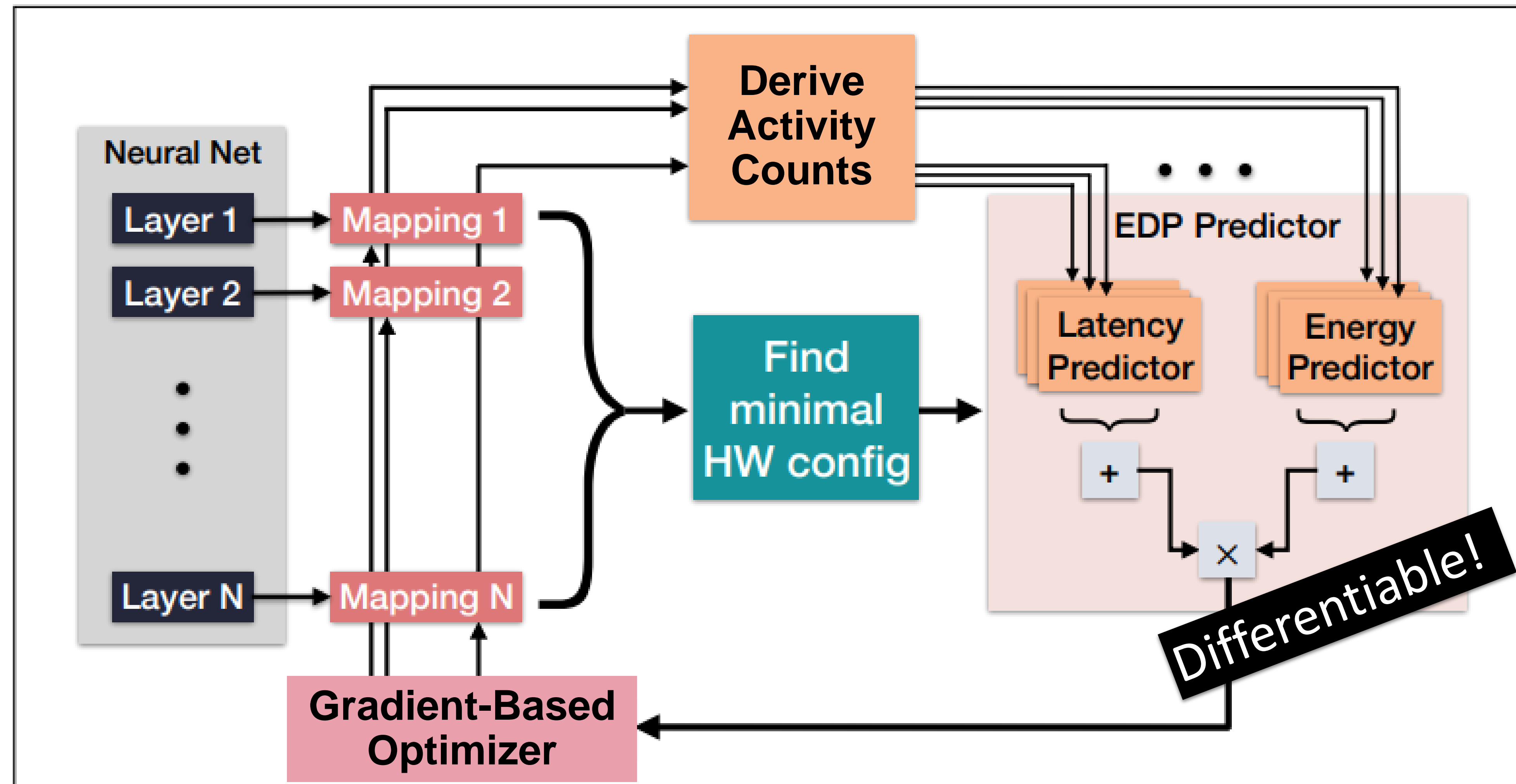
DOSA



- DOSA explicitly expresses the latency and energy performance as a **differentiable function** of mappings to enable gradient-based optimization.

# A DIFFERENTIABLE DSE FORMULATION

DOSA

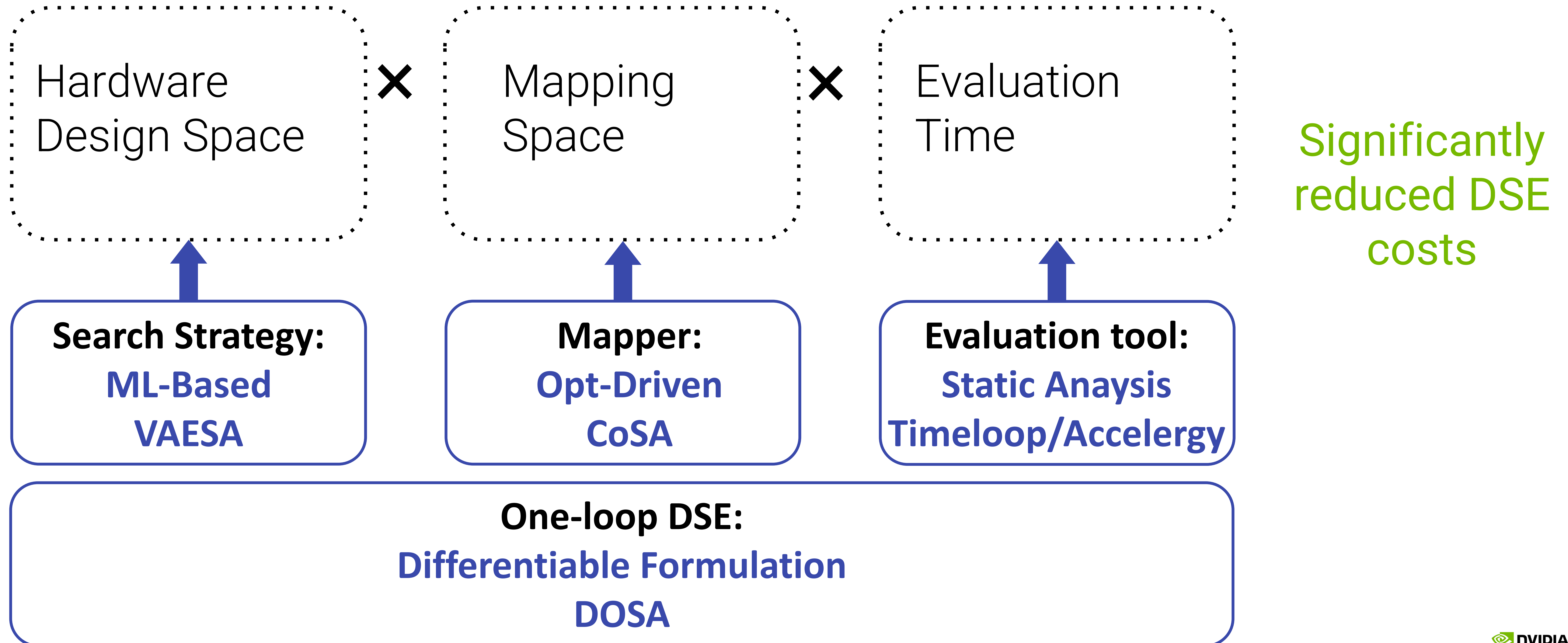


- Hardware design configurations are derived from the optimized in the one-loop mapping first DSE.
- DOSA's differentiable analytical model accurately predicts the performance and addresses the generalizability issues of the data-driven counter parts.
- It further improves the search efficiency and QoR in DSE.



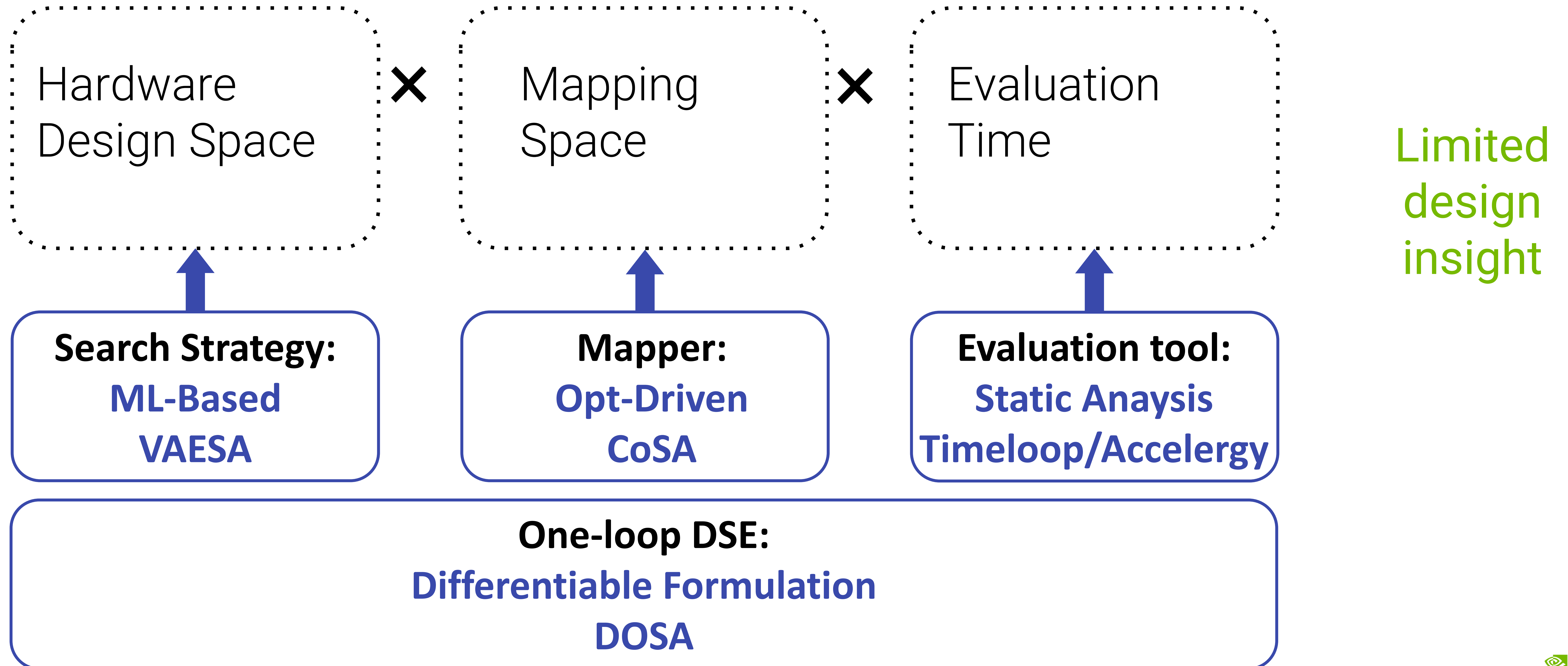
# DESIGN SPACE EXPLORATION TOOLS

Our approach



# DESIGN SPACE EXPLORATION TOOLS

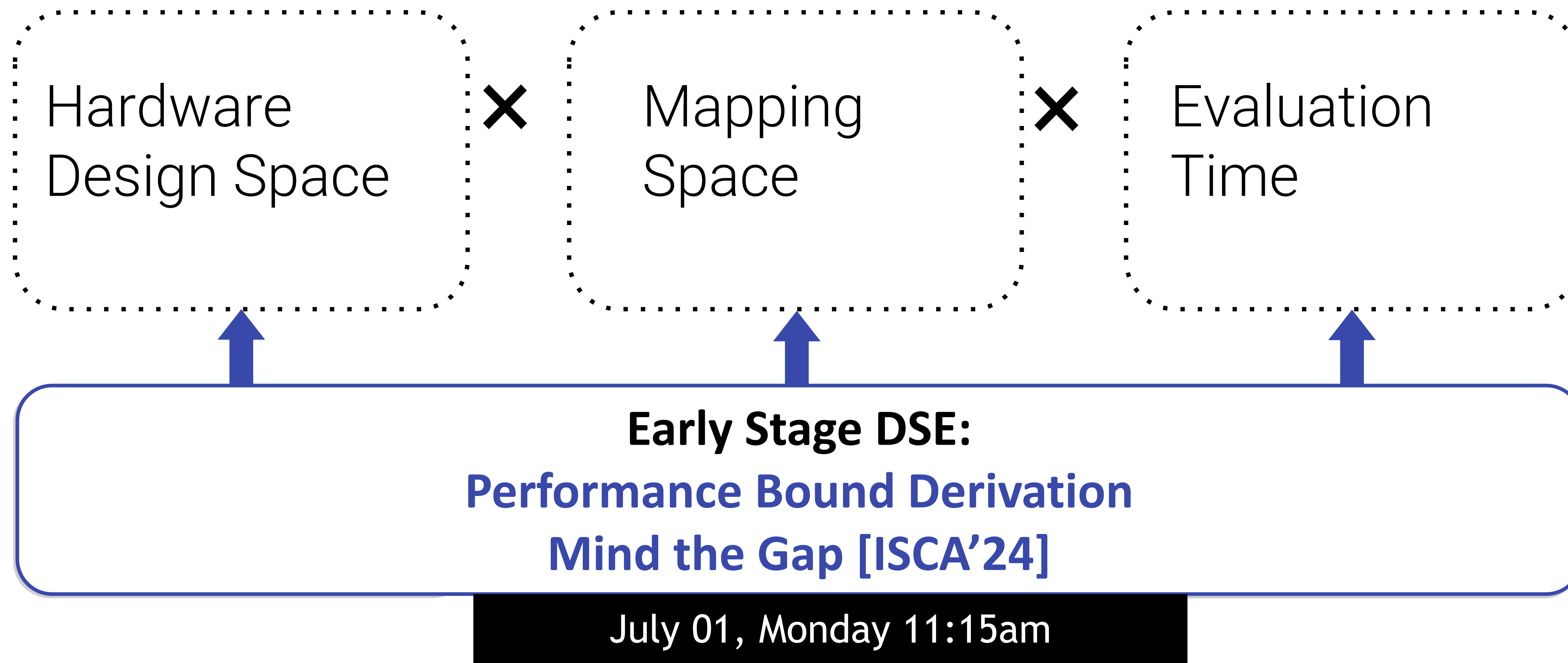
Our approach





# DESIGN SPACE EXPLORATION TOOLS

Our approach



# OPEN CHALLENGES AND OPPORTUNITIES

## #1 DSE for diverse workloads

- irregular
- input-dependent
- multi-tenant

## #2 DSE for dynamic system components

- SW/OS schedulers
- caching/paging

## #3 DSE for different execution models

- selection among cpu, vector, tensor units
- customization vs programmability tradeoffs

## #4 DSE generalizability for new workloads and constraints

- no consensus on DSE tools
- no uniform abstraction
- high-quality data is limited





THANK YOU

QIJING JENNY HUANG, NVIDIA

[jennyhuang@nvidia.com](mailto:jennyhuang@nvidia.com)