# Common barriers to deployment on edge devices

**Inability to deploy on edge devices**

**Unsatisfactory Accuracy or performance**

**Long development cycle**

deci.

# Models' power hunger is increasing rapidly
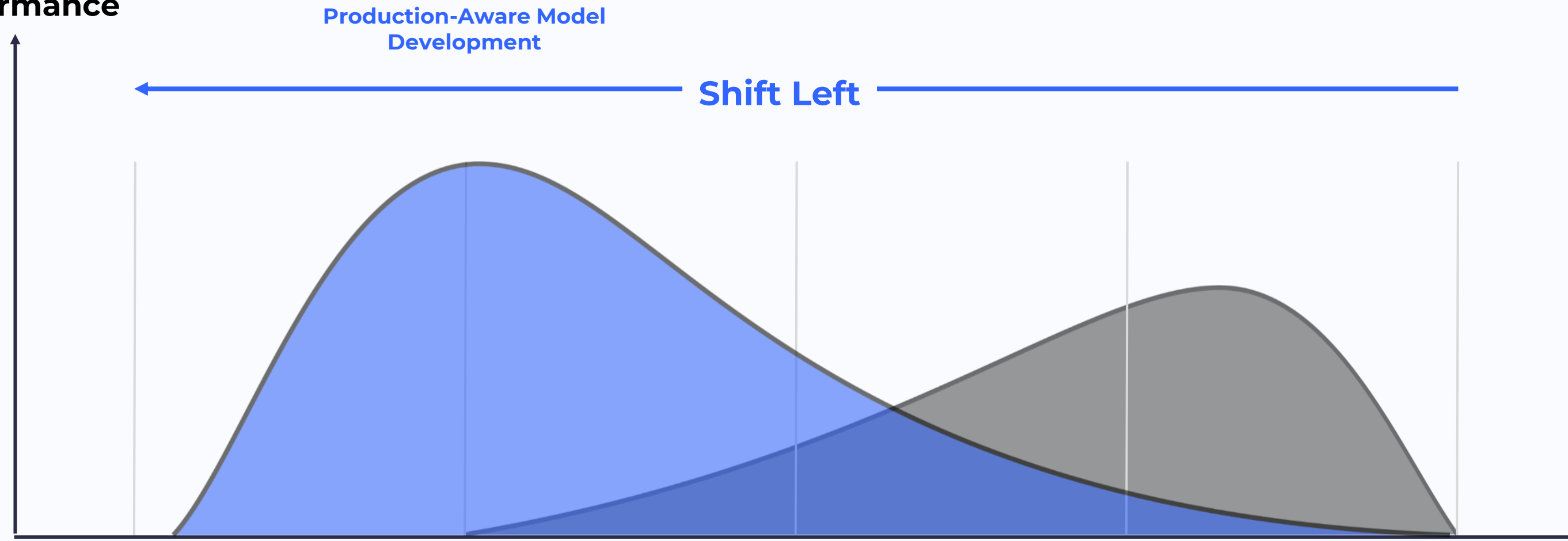
Model Complexity

Model Size
(# of FLOPs)

**The AI Efficiency Gap leads to:**

✗*Insufficient accuracy*
✗*High latency*
✗*Low throughput*
✗*Large model size*
✗*Large memory footprint*

Compute Power
(# of FLOPs / Sec)

Time

deci.

# AI Efficiency calls for a new development paradigm

**Attention to Inference Performance**

**Production-Aware Model Development**

**Shift Left**



Data Collection & Preparation

Model Selection & Development

Training

Optimization

Compilation Quantization Pruning

Deployment, Monitoring & Retraining

**Deep Learning Development Lifecycle**

Confidential

deci.

# Deci Deep Learning Development Platform

Powered by Neural Architecture Search

- ■ **Outperform SoTA with Custom NN Architectures**
  Save time and guarantee success by building accurate & fast architectures tailored
  for your performance targets & hardware

- ■ **Fast and Efficient Training Library**
  - Easily leverage advanced training techniques (Quantization Aware Training, Knowledge distillation)
  - Get SOTA hyperparameter recipes

- ■ **Automated Compilation & Quantization**
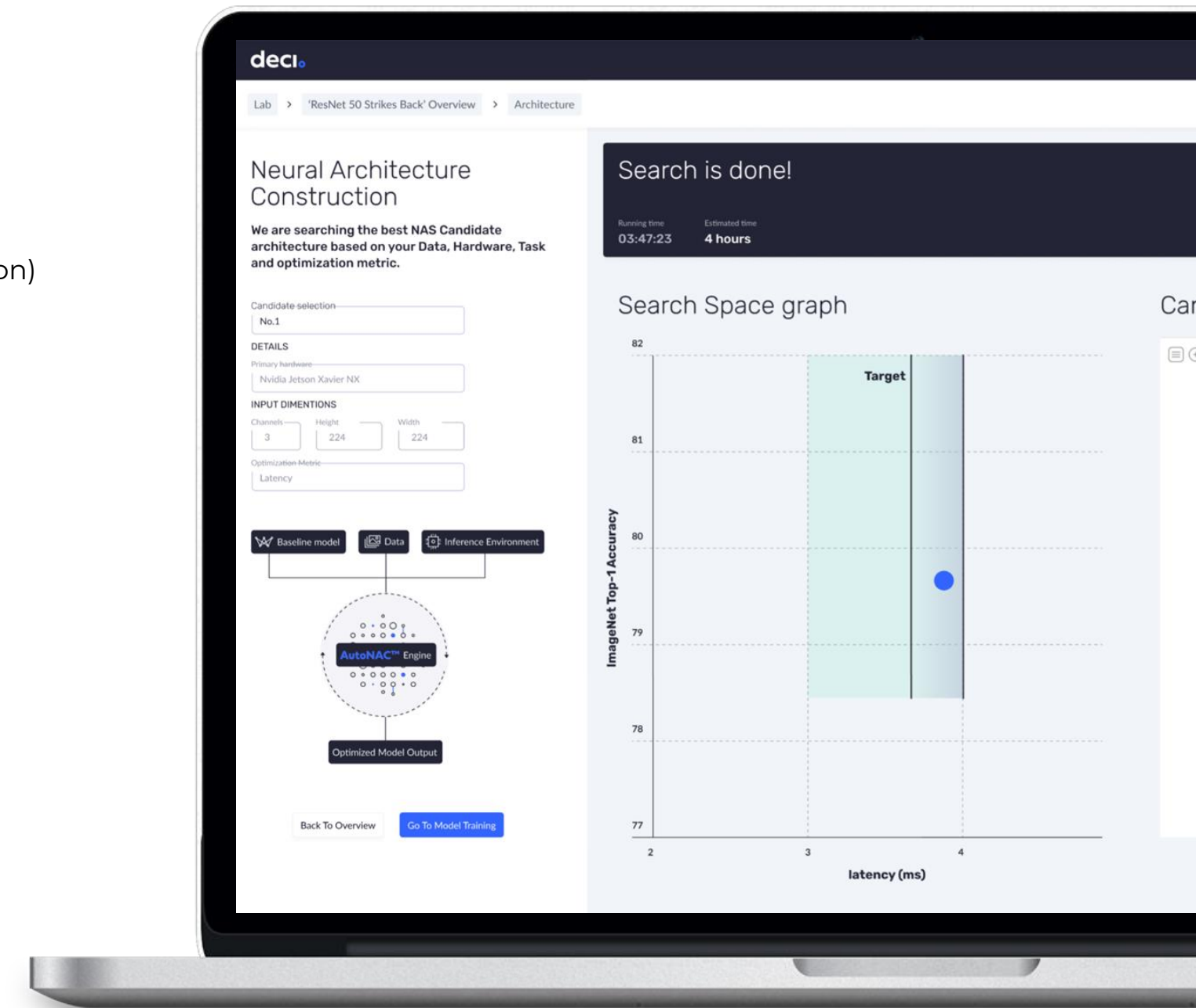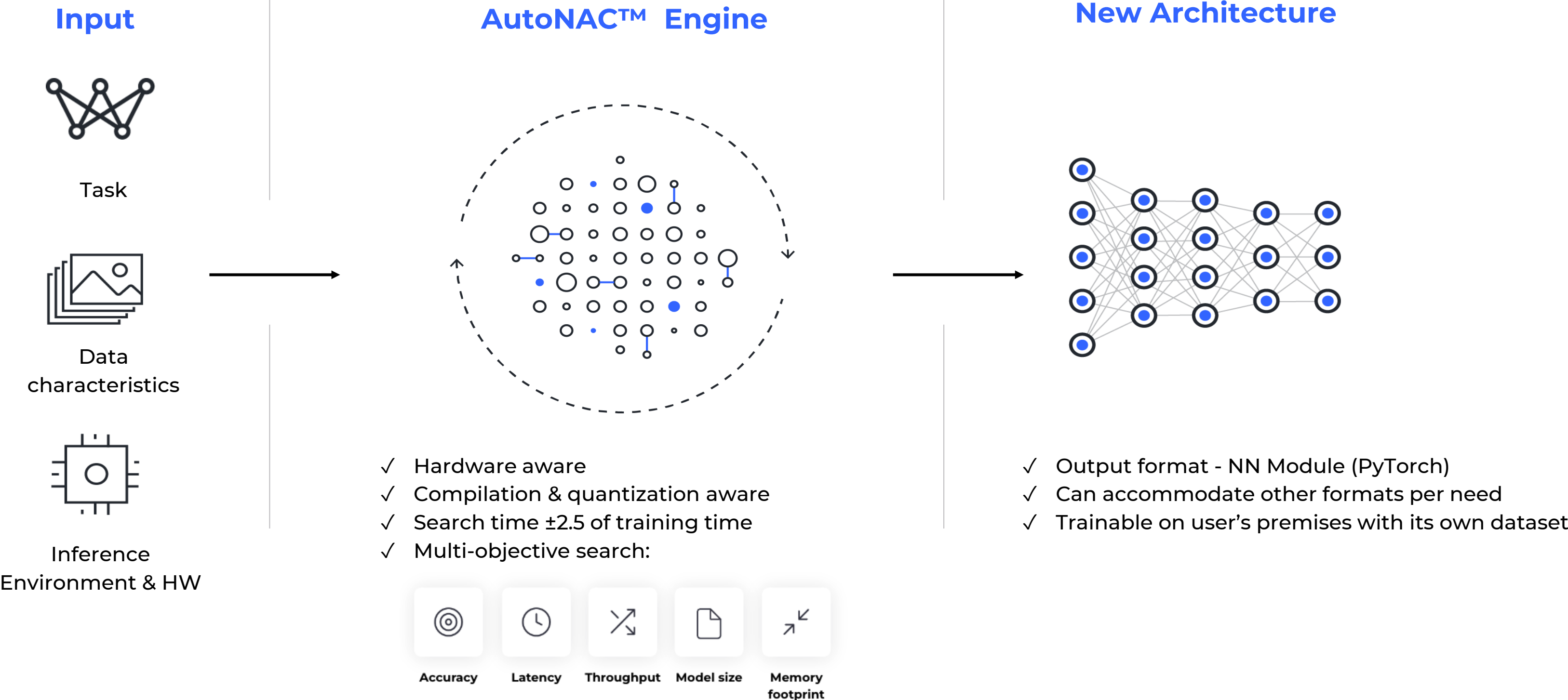  Optimize your trained models for your HW with a click of a button

- ■ **Inference Engine**
  Deploy with 3 lines of code using Deci's Python Inference Runtime Engine
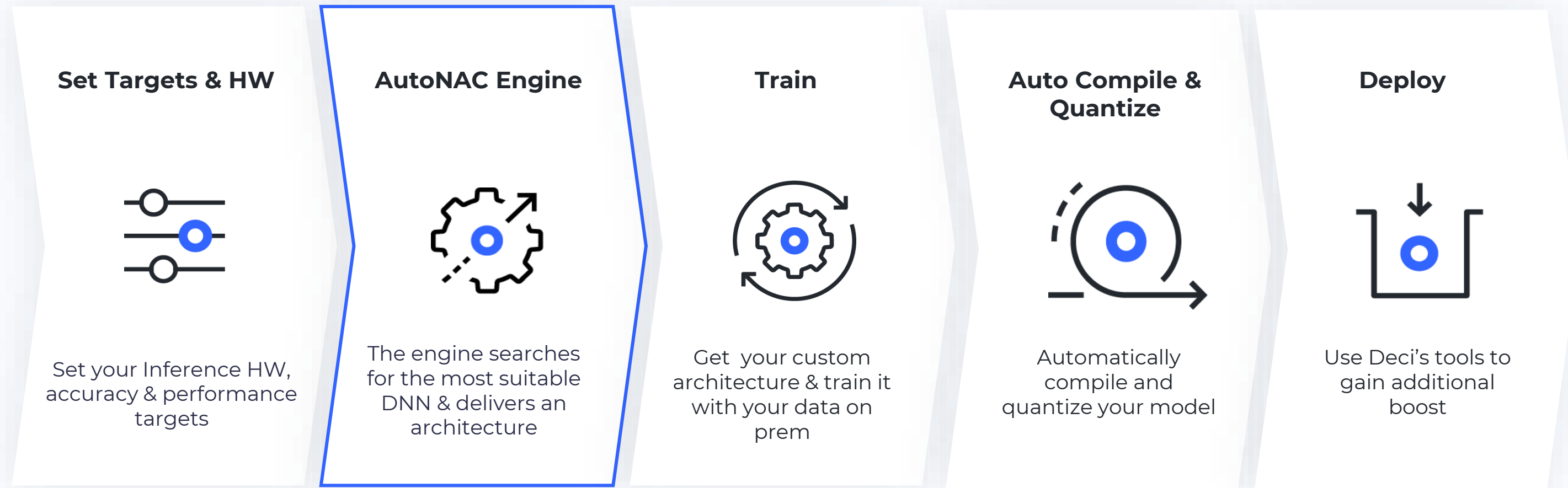
- ■ **Expert Support**
  Dedicated deep learning expert support

deci。

# Deci's AutoNAC Engine: Hardware-Aware Neural Architecture Search for DL Inference Efficiency

## Input

Task

Data characteristics

Inference Environment & HW

## AutoNAC™ Engine



✓ Hardware aware
✓ Compilation & quantization aware
✓ Search time ±2.5 of training time
✓ Multi-objective search:

Accuracy   Latency   Throughput   Model size   Memory footprint

## New Architecture

✓ Output format - NN Module (PyTorch)
✓ Can accommodate other formats per need
✓ Trainable on user's premises with its own dataset

deci.

# Build custom models with Deci

### Set Targets & HW

Set your Inference HW, accuracy & performance targets

### AutoNAC Engine

The engine searches for the most suitable DNN & delivers an architecture

### Train

Get your custom architecture & train it with your data on prem

### Auto Compile & Quantize

Automatically compile and quantize your model

### Deploy

Use Deci's tools to gain additional boost

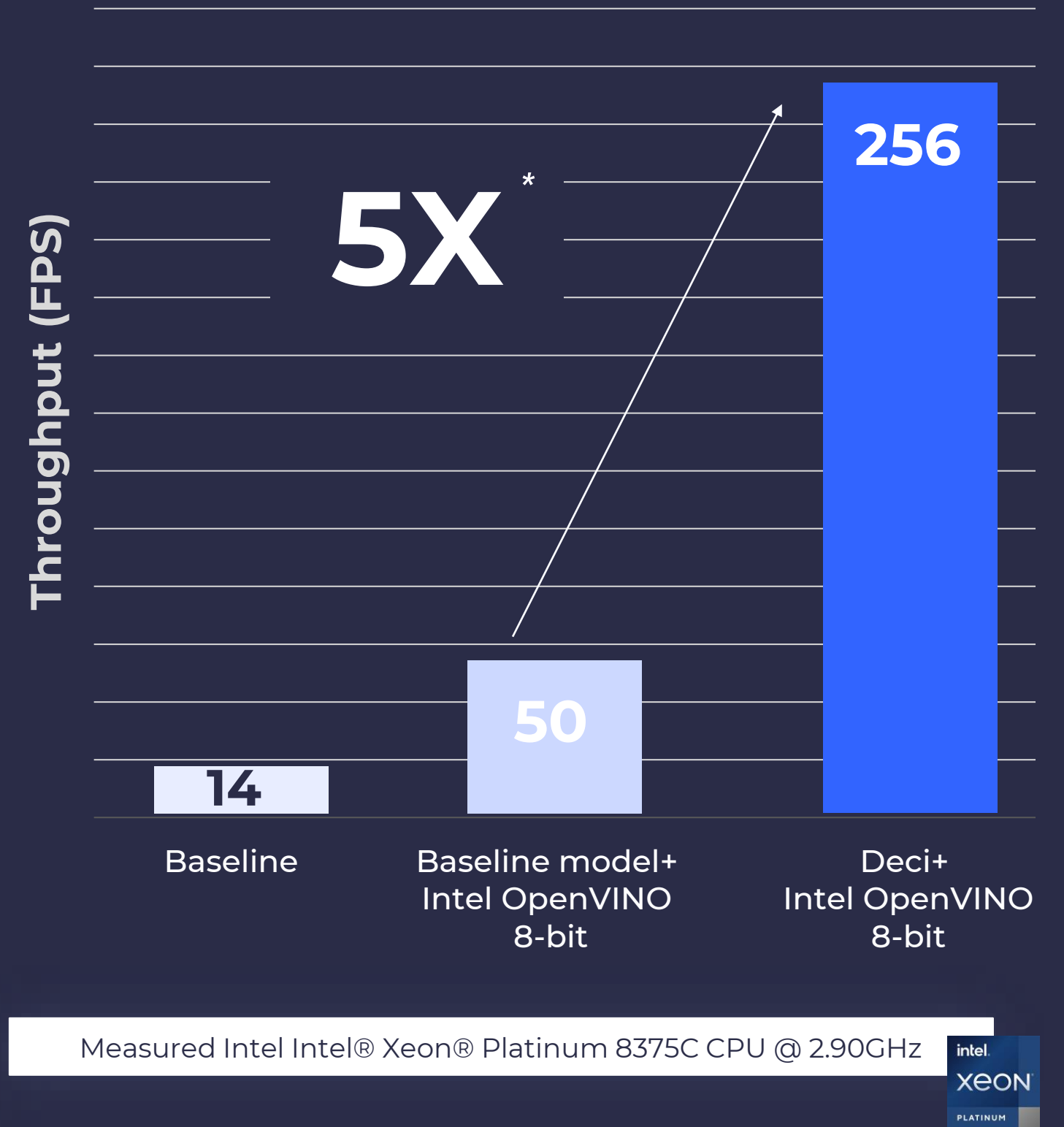**Performance Targets**

Accuracy Preserving

Latency

Throughput

Model size

Memory footprint

deci.

# An **inherent advantage** driven by Deci's algorithmic level optimization

**MLPerf**

**Hardware Aware Model Architecture construction** ———— Neural Architecture Search

**Model Compression Techniques** ———— Quantization / Pruning

**Runtime Optimization** ———— Graph Compilers / Drivers

**Inference Hardware** ———— FPGA / GPU / ASIC / Embedded / CPU / TPU

## A comparison of ResNet 50 & DeciNet on ImageNet

**5X** *

**Throughput (FPS)**

| | |
|---|---|
| Baseline | **14** |
| Baseline model+ Intel OpenVINO 8-bit | **50** |
| Deci+ Intel OpenVINO 8-bit | **256** |

Measured Intel Intel® Xeon® Platinum 8375C CPU @ 2.90GHz

intel XEON PLATINUM

*Performance claim data: https://community.intel.com/t5/Blogs/Tech-Innovation/Artificial-Intelligence-AI/Intel-Collaboration-With-Deci-Boosts-AI-Performance-on-Intel/post/1404029

deci.

# With **Deci**, You can Build **Better** Models, **Faster**.

## Gain Unparalleled Inference Performance

Up to **5X** *

acceleration

## Shorten Time to Market

**3** weeks

on average to reach a production-ready model

## Guarantee Success In Production

**Built for purpose & Expert Support**

*Performance claim data: https://community.intel.com/t5/Blogs/Tech-Innovation/Artificial-Intelligence-AI/Intel-Collaboration-With-Deci-Boosts-AI-Performance-on-Intel/post/1404029

deci.

# Use Cases - How AI teams are using Deci?

## Enables Inference on Edge Devices

Enable inference on resource constrained devices ( e.g. Edge devices, mobile etc.)

## Boost Inference Performance

Outperform SOTA models with better accuracy, latency, throughput, smaller memory footprint & model size.

## Reduce Training & Inference Costs

Maximize Hardware utilization. Make the of most of your current hardware or more to a more affordable one. Cut up to 80% of your cloud costs.

## Simplify Development, Shorten Time to Market

Automate model development & optimization steps. Eliminate uncertainty, guarantee success in production and reach production faster.

deci.

# Thank You.

deci.

# Notices and Disclaimers