# Empower Your Edge Device Using NetsPresso—No AI Engineer Required!

Tae-Ho Kim

CTO/Co-founder

Nota AI

# NetsPresso

Nota AI

# Efficient AI = Compressed AI

- AI compression is essential for commercialization

**70MB**

**Recent AI model size**

**350KB**

**Embedded HW capacity**

*"I think one of the <u>hardest technical problems</u> …*
*<u>the neural net needs</u>*
*<u>to be</u> <u>**compressed**</u> <u>into a fairly small computer</u>,*
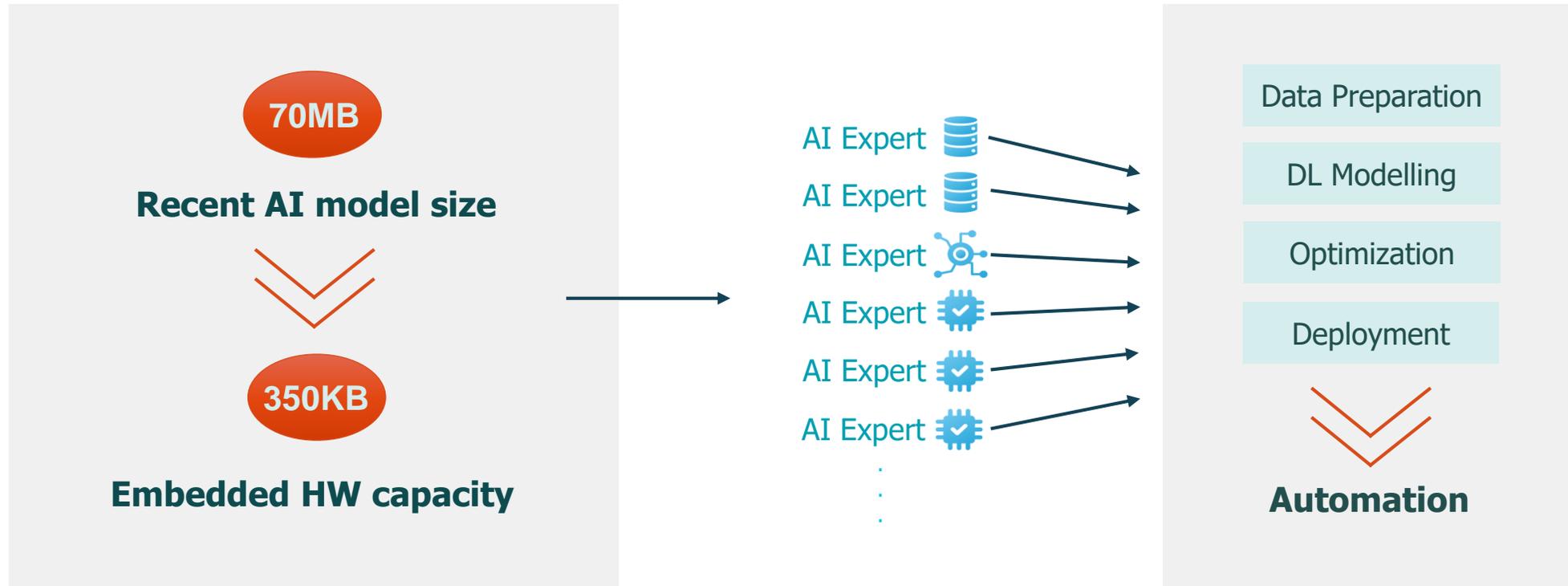*a very efficient computer that was designed, …"*

By Elon Musk, CEO at Tesla / <Tesla 2021 1Q Earning Call>
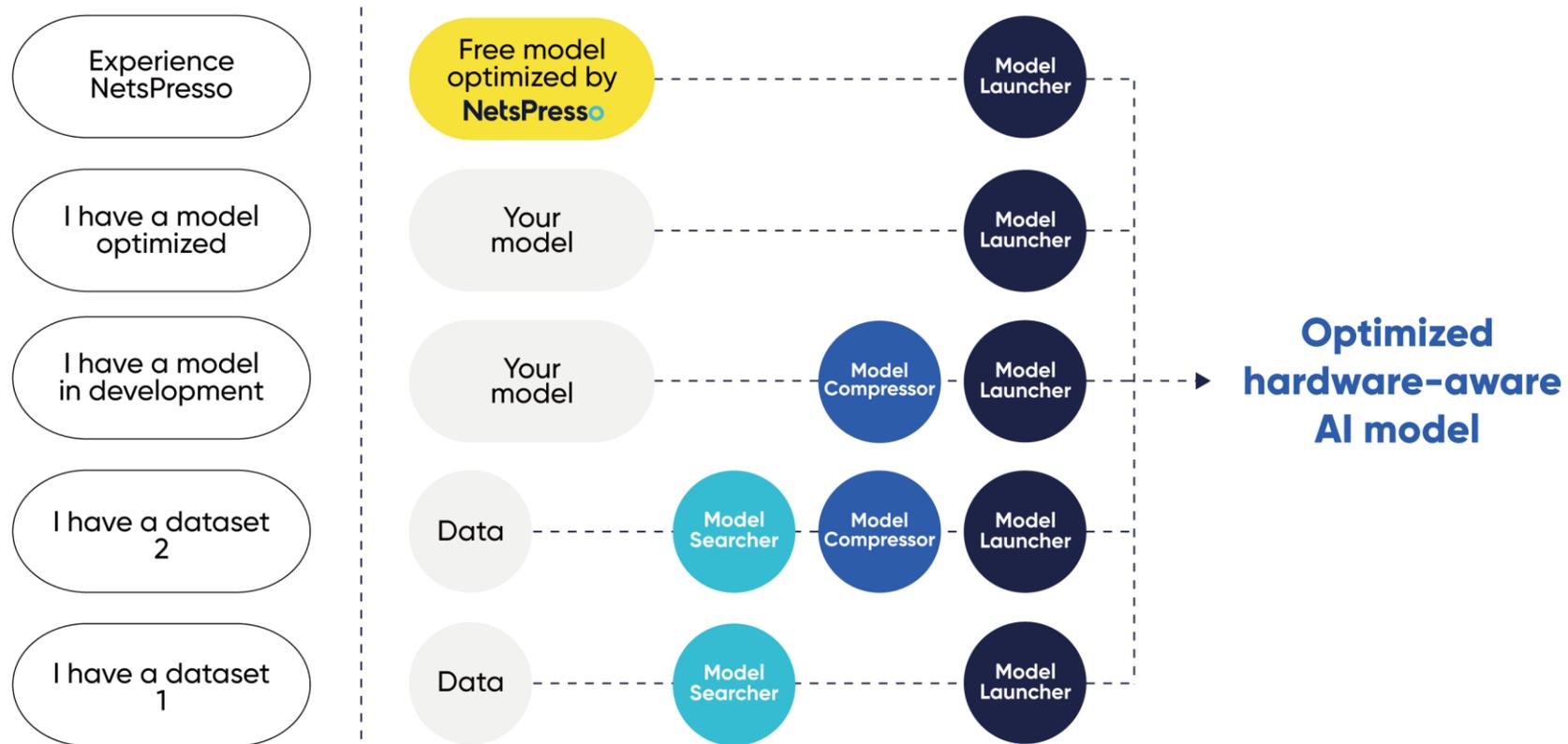
# Efficient AI = Efficient Automation

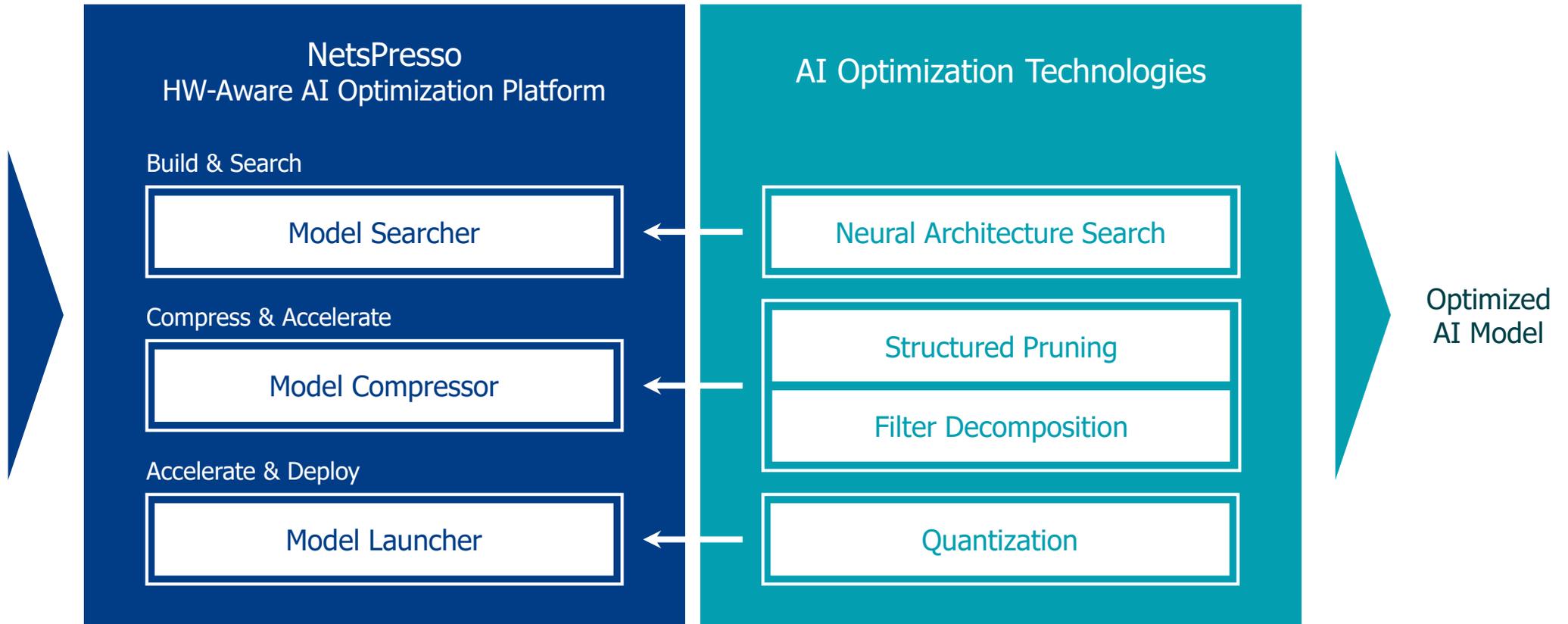- AI expert replacement is the essential for future AI



*Replacement*

70MB

**Recent AI model size**

350KB

**Embedded HW capacity**

AI Expert
AI Expert
AI Expert
AI Expert
AI Expert
AI Expert

Data Preparation

DL Modelling

Optimization

Deployment

**Automation**

# NetsPresso: AI Optimization Platform

- Which Module do you need?

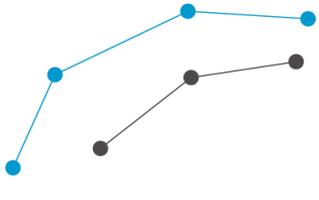# NetsPresso and AI Optimization Technologies



**NetsPresso**
HW-Aware AI Optimization Platform

Build & Search
- Model Searcher

Compress & Accelerate
- Model Compressor

Accelerate & Deploy
- Model Launcher

**AI Optimization Technologies**
- Neural Architecture Search
- Structured Pruning
- Filter Decomposition
- Quantization

Target Performance
Target Hardware

Optimized
AI Model

# Overview of Each Module



**Model Searcher**
- Optimal trade-off models
- 5X faster model
- Optimization targets: NVIDIA, arm, Raspberry Pi, intel, Coral

**Model Compressor**
- Model profiling
- Compression methods: Pruning, Filter decomposition
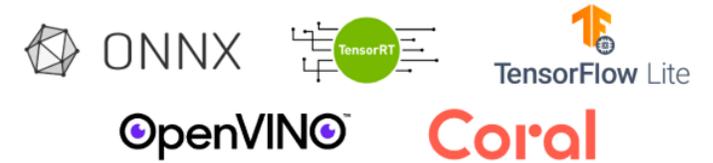- 30X compressed

**Model Launcher**
- Ready to deploy: Pre-processing, Inference engine, Post-processing
- Benchmarks on device: Latency, Power, Memory consumption
- Model converting: ONNX, TensorRT, TensorFlow Lite, OpenVINO, Coral

## Hardware-Aware AutoML

- NetsPresso Model Searcher automatically searches optimized models for your target device

**Benefits**
- Reducing development time from months to weeks
- Provides multiple models with various options
- Provides model close to production level based on actual HW test
- Creates a model with low latency

Dataset

Model Searcher

Optimal trade-off models

5X faster model

Optimization targets

NVIDIA    arm

intel    Raspberry Pi

Optimized AI Model

# Module 1: NetsPresso Model Searcher



Dataset: MS COCO

Target HW: NVIDIA Jetson Nano

- Comparison with other popular models
- The closer it is to the top left, the better the result
- Model Searcher offers several options to make the most suitable choice on the trade-off

## Ready-to-use toolkit

- Make the compression process easy and fast

**Benefits**
- Supports all CNN architectures
- Optimal compression ratio is recommended
- Eliminates months of paper implementation time
- Minimal loss of information

AI Model

**Model Compressor**

Model profiling

Compression methods

Pruning

Filter decomposition

30X compressed

Compressed Model

# Module 2: NetsPresso Model Compressor

- ## Best practices

### Classification

| Model | Method | Accuracy (%) | FLOPs (M) | Params (M) |
|---|---|---|---|---|
| ResNet50 | Original | 78.03 | 2596.06 | 23.71 |
| | *NPTK | 76.63 (-1.4) | 224.70 (11.55x) | 2.17 (10.91x) |

### Object Detection

| Model | Method | mAP(0.5) (%) | FLOPs (M) | Params (M) |
|---|---|---|---|---|
| YOLOv4 | Original | 82.22 | 61871.82 | 262.90 |
| | *NPTK | 87.23 (+5.01) | 11459.69 (5.4x) | 2.75 (7.49x) |

### Super Resolution

| Model | Method | PSNR (dB) | FLOPs (M) | Params (M) |
|---|---|---|---|---|
| EDSR | Original | 31.95 | 228665.89 | 1.52 |
| | *NPTK | 31.61 (-0.34) | 83124.98 (2.75x) | 0.72 (2.1x) |

*NPTK: NetsPresso Model Compressor

NotaAI

© 2022 Nota AI

## Convert and Package for the Deployment

- Benchmark the model on device immediately
- Deploy the model on device immediately

**Benefits**
- Ease of model converting
- Performance benchmarks and recommendations on actual devices
- Ready-to-deploy packaging

Optimized AI Model

**Model Launcher**

Ready to deploy
- Pre-processing
- Inference engine
- Post-processing

Benchmarks on device
Latency Power Memory consumption

Model converting
ONNX  TensorRT  TensorFlow Lite
OpenVINO  Coral

Production-Ready Package

**Deep Learning Engineer**

Trained model →

← HW benchmarks (latency, power consumption)

← Ready-to-deploy package

Convert and accelerate the model on *actual* devices

**Model converting**

ONNX    TensorRT    TensorFlow Lite
OpenVINO    Coral

**Optimization targets**

NVIDIA    arm    Raspberry Pi
Coral    intel

Ready-to-deploy packaging

NotaAI

# Application Solutions

# Application Solutions Powered by NetsPresso

**HW-aware MLOps Platform**

**Intelligent Transportation System**

**Driver Monitoring System**

# Intelligent Transportation System (ITS)

## AI Camera



- Collection of traffic data including car type, traffic flow robust to weather conditions
- Collection of spatial information including queue length, # of queueing vehicles, occupancy, average speed etc.

## AI Traffic Signal Control



- Reinforcement Learning-based optimal control using traffic data from AI Traffic Camera
- Isolated intersection control
- Arterial control
- Traffic network control

## AI Safe Crossing



- Building LDM (Local Dynamic Map) layer 4 using short-term object information (< second)
- Prediction of potential accident/incident
- Providing warning signs (3-seconds before) through in-ground devices, such as light or voice alerts.

# ITS Use Case: AI Traffic Camera

# Driver Monitoring System (DMS)

- DMS detects unsafe drivers' behaviors
- We provide each feature as an independent API
- Modules can be chosen and packaged by request



**Drowsiness**

**Distraction**

**Unregistered Driver**

**Using Cell Phone**

**Smoking**

18

# DMS Use Case: Solution Applied to a Low-spec Device

| Facial Recognition | Face Detection | Landmark Detection |
| --- | --- | --- |
| Behavior Analysis | Video Enhancement | Pose Estimation |

**NetsPresso**

**Nota AI**

# Solution Suggestions

If you <u>have only a dataset</u> and no AI experts ▶ **Model Searcher**

If you <u>have the model</u> but need optimization ▶ **Model Compressor**

If you need <u>to optimize</u> on your HW ▶ **Model Launcher**

If you <u>need</u> the AI application ▶ **ITS**   **DMS**

# More Information

## Website

- http://nota.ai

- http://netspresso.ai

## Oral Session

📅 **Tuesday, May 17**

1:30 pm

Empower Your Edge Device Using NetsPresso – No AI Engineer Required!

## Booth

# #418

# Thank You!

For further inquires, please contact

Tae-Ho Kim, thkim@nota.ai

Visit www.nota.ai for detailed information

Nota AI

# Interface: NetsPresso Model Searcher

# Interface: NetsPresso Model Searcher

# Interface: NetsPresso Model Searcher

netspresso

- Model Search
- Datasets
- Projects
- Models

- Documentation
- Github Discussion

**Target device**

**Target device ***
- ◉ NVIDIA Jetson — AGX Xavier ▾
- ○ Raspberry Pi
- ○ Intel Xeon

**Output Format**

**Framework ***
Tensor RT 8.0.1 ▾

**SW version ***
JetPack 4.6 ▾

**Output datatype ***
- ◉ FP32  ○ FP16  ○ INT8  ○ INT4

**Inference batch size ***

ⓘ • Support range: 1~32
• TFlite only supports batch size 1

**Model training**

**Target latency (ms)***

# Interface: NetsPresso Model Compressor

# Interface: NetsPresso Model Compressor

# Interface: NetsPresso Model Compressor

# Interface: NetsPresso Model Launcher

# Interface: NetsPresso Model Launcher