# ACCELERATE YOUR INFERENCING WITH INTEL® DEEP LEARNING BOOST

**Shailen Sobhee**

**AI Software Technical Consultant**

shailen.sobhee@intel.com

# Audience pre-requisites

- Familiar with deep learning stages
    - Training and inferencing
- Have a basic knowledge about hardware
    - know what are vector registers, like AVX-512

# Outline

- What is Intel® Deep Learning Boost (Intel® DL Boost)
- Why is Intel® DL Boost useful?
- What are **V**ector **N**eural **N**etwork **I**nstructions (**VNNI**)
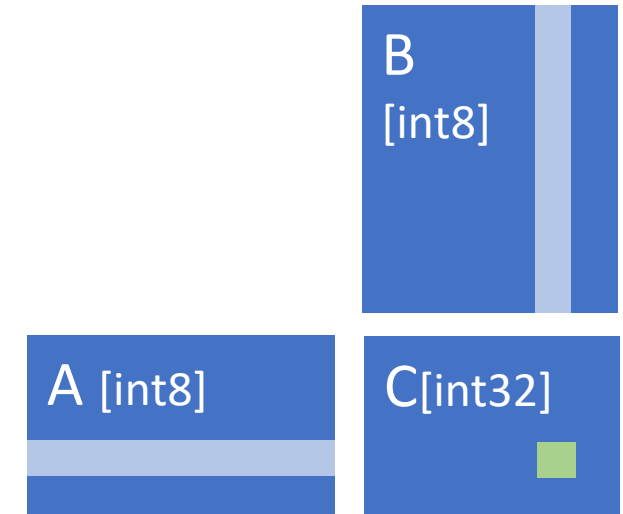- Sample results

# What is Intel® Deep Learning Boost ?

Intel® DL Boost:

- extends the AVX-512 instructions

- designed to deliver **significant** and **more efficient** Deep Learning (Inference) acceleration

- for deep learning workloads optimized to use the Vector Neural Network Instruction (VNNI)

- on Intel® Xeon® Scalable processor

- as from the 2$^{nd}$ generation (codename "Cascade Lake")
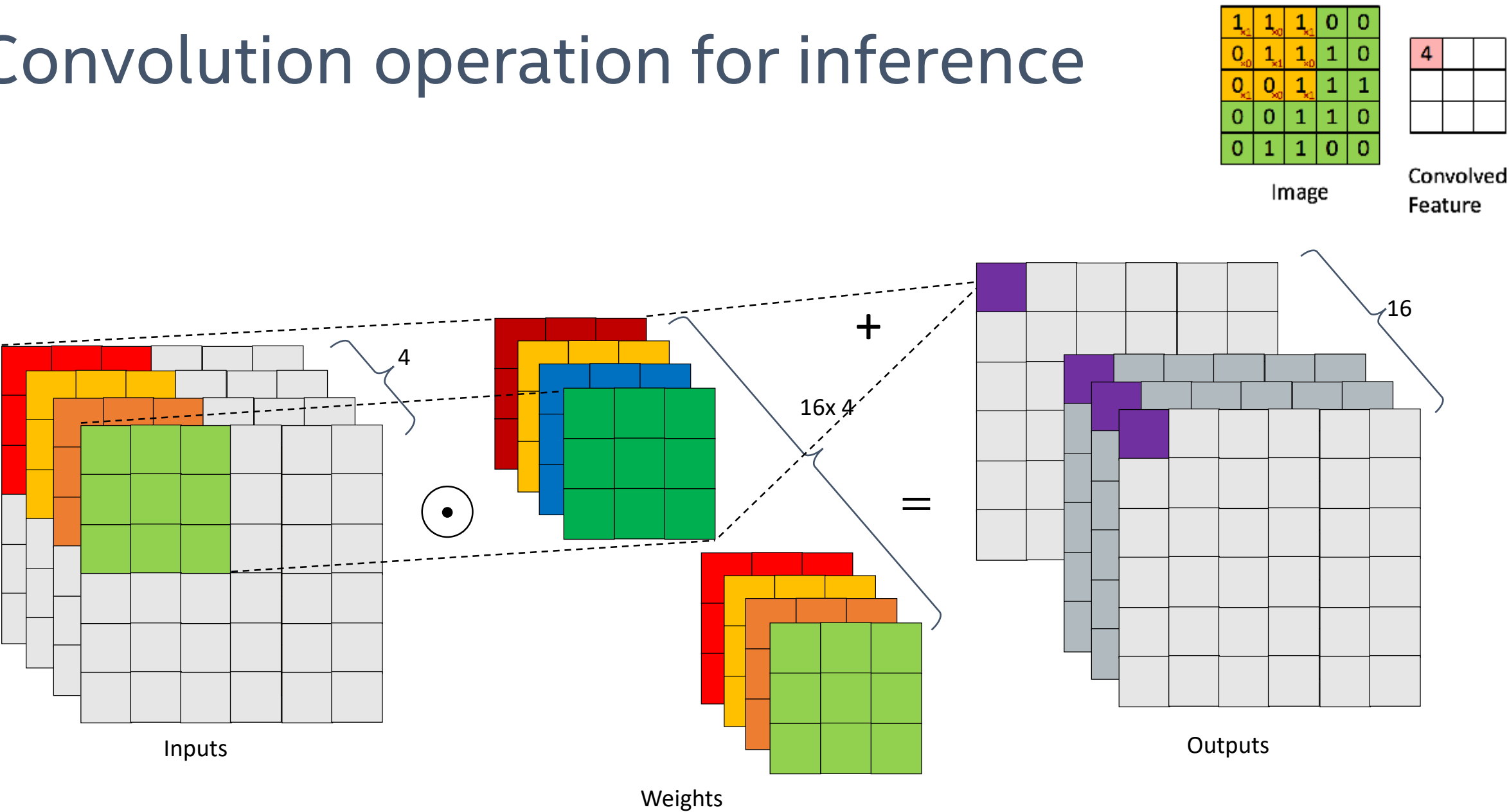
# Deep Learning Foundations

- Heavy compute (Matrix Multiplications) are the foundation of many DL applications
    - **Multiply** a row*column values, **accumulate** into a single value
- Traditional HPC and many AI training workloads use floating point
    - Massive dynamic range of values (FP32 goes up to ~2^128)

B
[int8]

A [int8]

C[int32]

**Matrix Multiply**

**A x B = C**

# Convolution operation for inference



Inputs ⊙ Weights = Outputs

4    16x 4    16

Image          Convolved Feature

# Why do we need Intel® Deep Learning Boost?

# The key term:

## Quantization

# Here's why Quantization matters

**Floating Point**
96.1924

**Integer**
96

32 -bit

8 bit

| | |
|---|---|
| 10110110 | 10110110 |
| 10110110 | 10110110 |

10110110

# Here's why Quantization matters

**Lower Power**

**Lower memory bandwidth**

**Lower storage**

**Higher performance**

**Important: Acceptable accuracy loss**

Image credits: see backup

# VNNI INSTRUCTION SET

Image

Convolved
Feature

# Here's one tool in your arsenal to do it ☺

# Intel® Distribution of OpenVINO™ in a nutshell

**1**  **2**  **3**  **4**

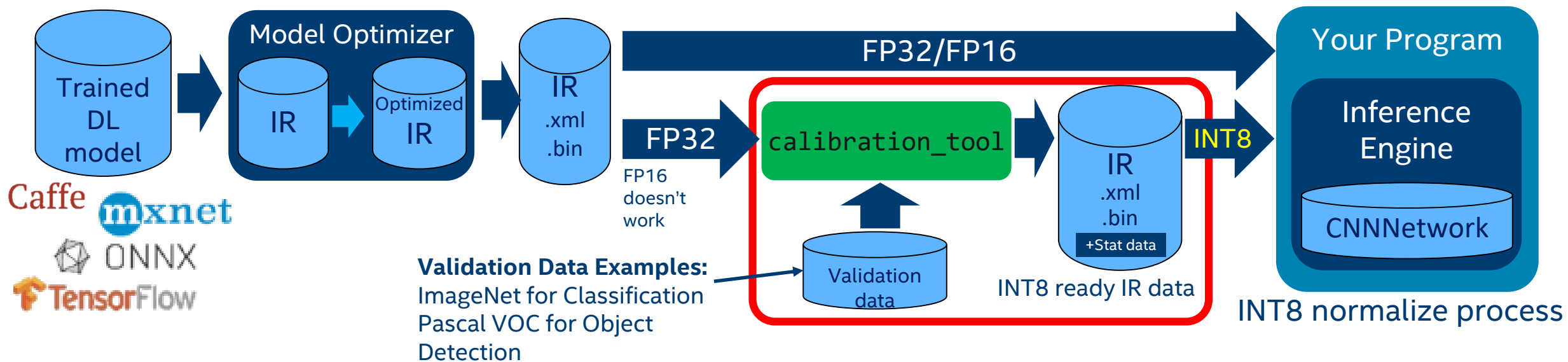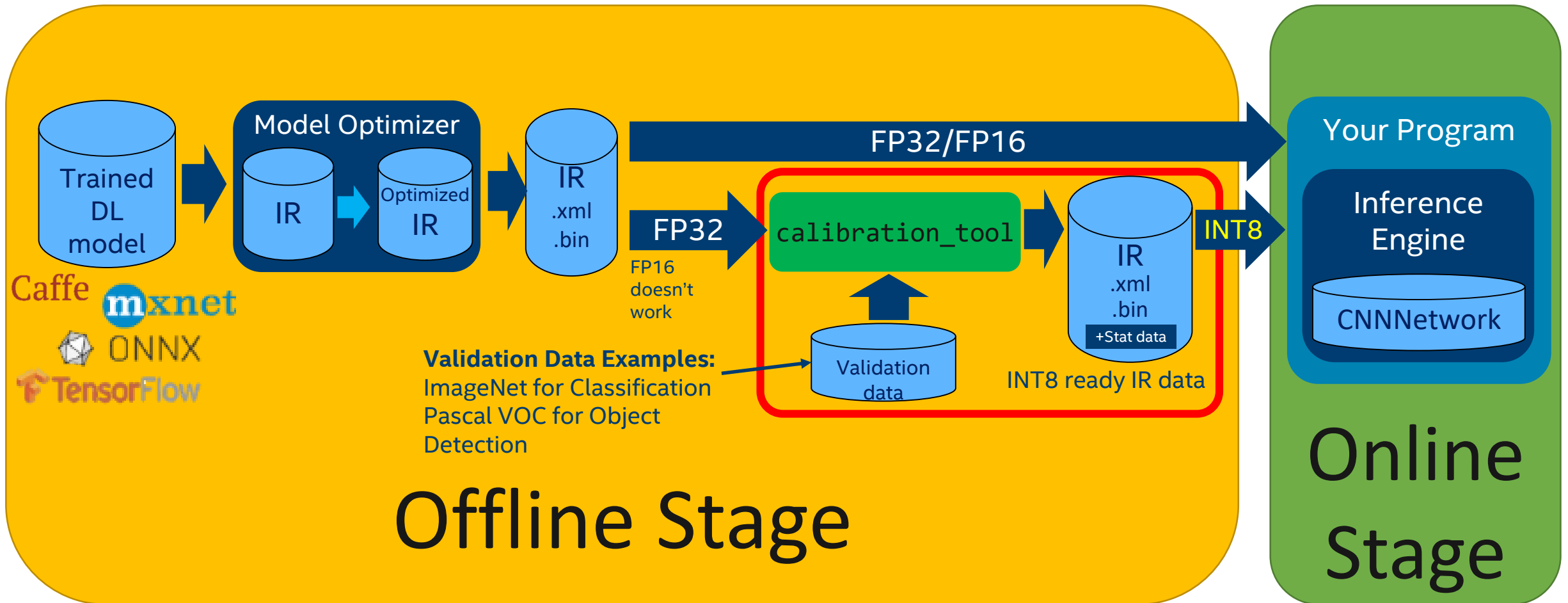# Intel® Distribution of OpenVINO™ in a nutshell

# Intel® Distribution of OpenVINO™ in a nutshell
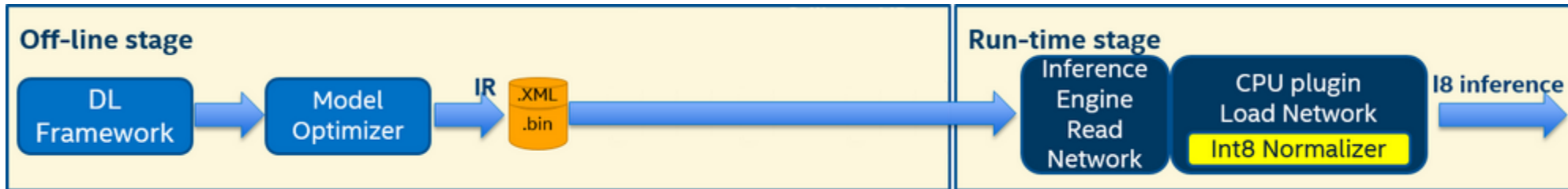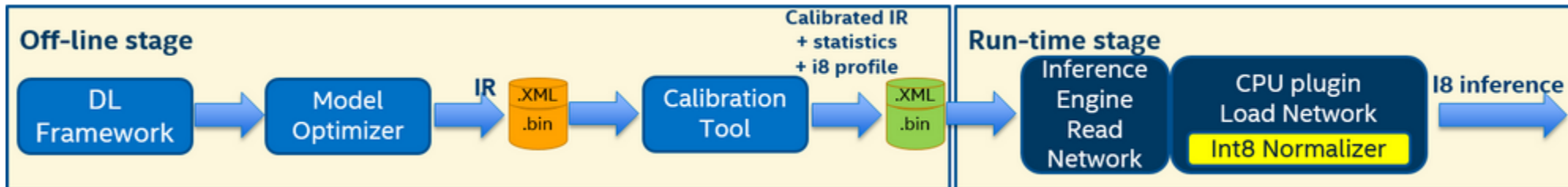
# Intel® Distribution of OpenVINO™ in a nutshell



https://docs.openvinotoolkit.org/latest/_docs_IE_DG_Int8Inference.html

# Sample results

## Demo 1



## Demo 2



Both executions on Intel® Cascade Lake CPU

# Sample results

```
FP32 Inference: 354.1443868062011
Int8 Inference: 798.4037049108333
Speed Up: 2.2544581663742274
```



FP32 vs Int8 Inference Speed

# Key take away

Try the **Intel® Distribution of OpenVINO™**: https://software.intel.com/en-us/openvino-toolkit

Benefit from **faster** inference speeds with **INT8** leveraging **VNNI** instructions on **Intel® Cascade Lake** CPUs.

# Summary

- What is Intel® Deep Learning Boost (Intel® DL Boost)
- What are **V**ector **N**eural **N**etwork **I**nstructions (**VNNI**)
- Why is Intel® DL Boost useful?
- Intel® Distribution of OpenVINO™