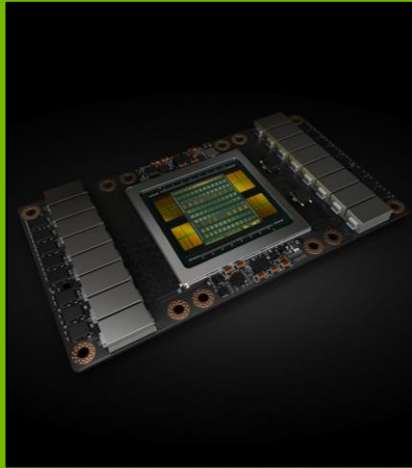# NVIDIA FOR DEEP LEARNING

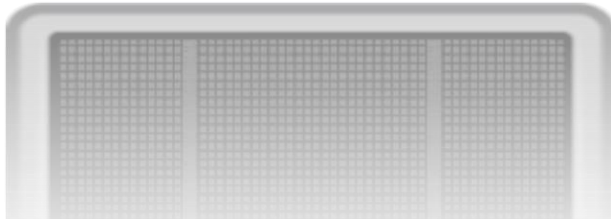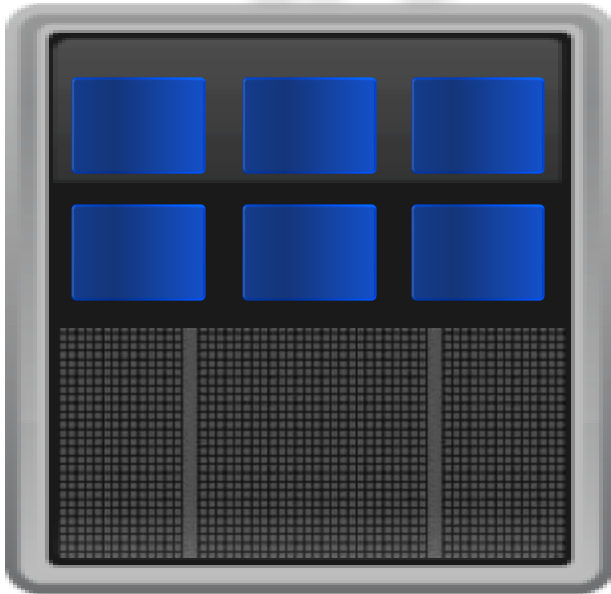Bill Veenhuis | bveenhuis@nvidia.com
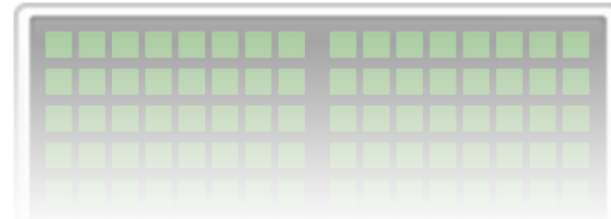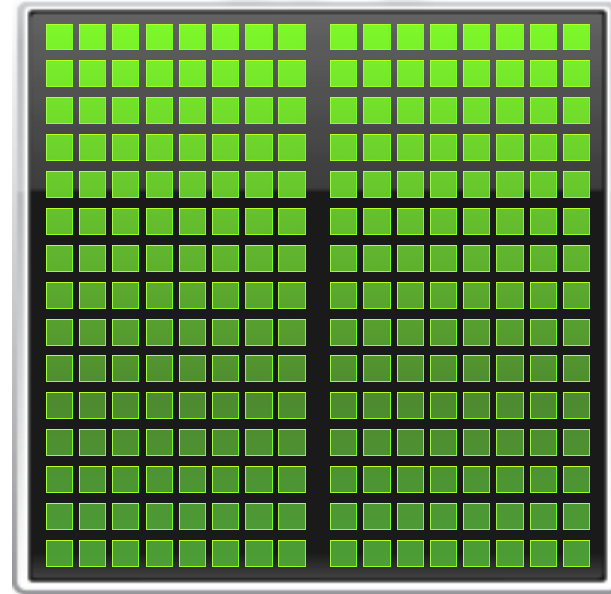
# Nvidia is the
# world's leading ai platform



ONE ARCHITECTURE — CUDA

# GPU: Perfect Companion
# for Accelerating Apps & A.I.

# AGENDA
# &
# TOPICS

- Intro to AI

- Deep Learning Intro

- NVIDIA's DIGITS

- Autoencoding Enhancement

- TensorRT

# Intro to AI

# ARTIFICIAL NEURONS

Biological neuron

Artificial neuron

impulses carried toward cell body

branches of axon

dendrites

nucleus

axon

axon terminals

impulses carried away from cell body

cell body

From Stanford cs231n lecture notes

$w_1$  $w_2$  $w_3$

$x_1$  $x_2$  $x_3$

$y$

Weights ($W_n$)
= parameters

$y = F(w_1 x_1 + w_2 x_2 + w_3 x_3)$

# ARTIFICIAL NEURAL NETWORK

A collection of simple, trainable mathematical units that collectively learn complex functions

Hidden layers

Input layer

Output layer

Given sufficient training data an artificial neural network can approximate very complex functions mapping raw data to output decisions
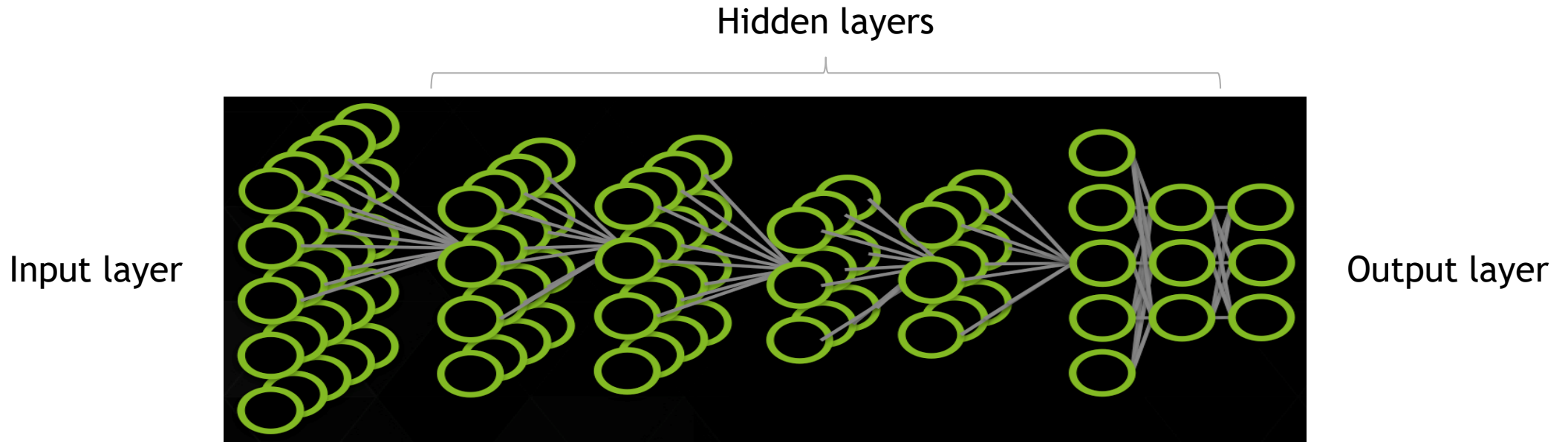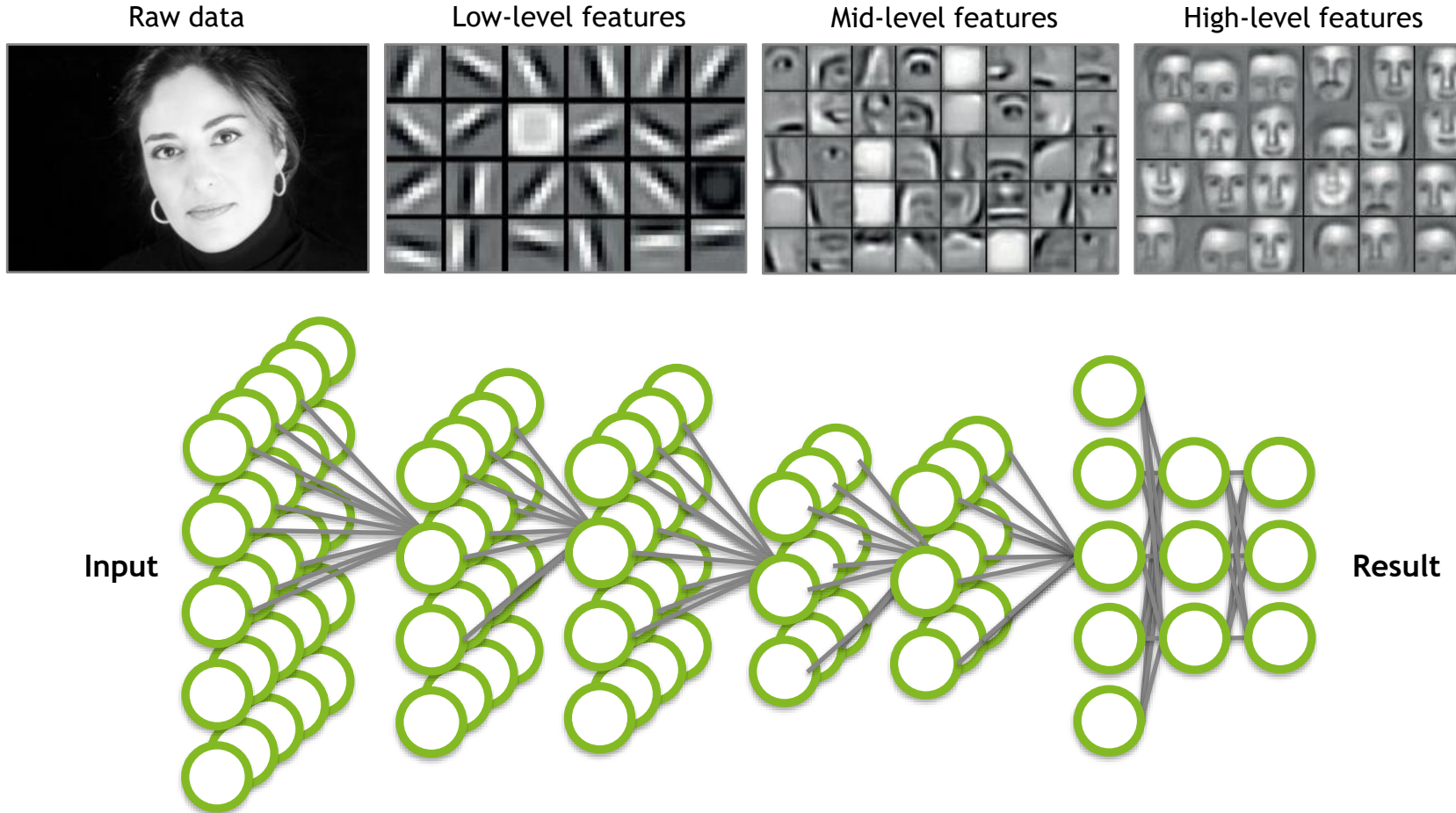
# DEEP NEURAL NETWORK (DNN)

Raw data     Low-level features     Mid-level features     High-level features



Input                     Result

**Application components:**

**Task objective**
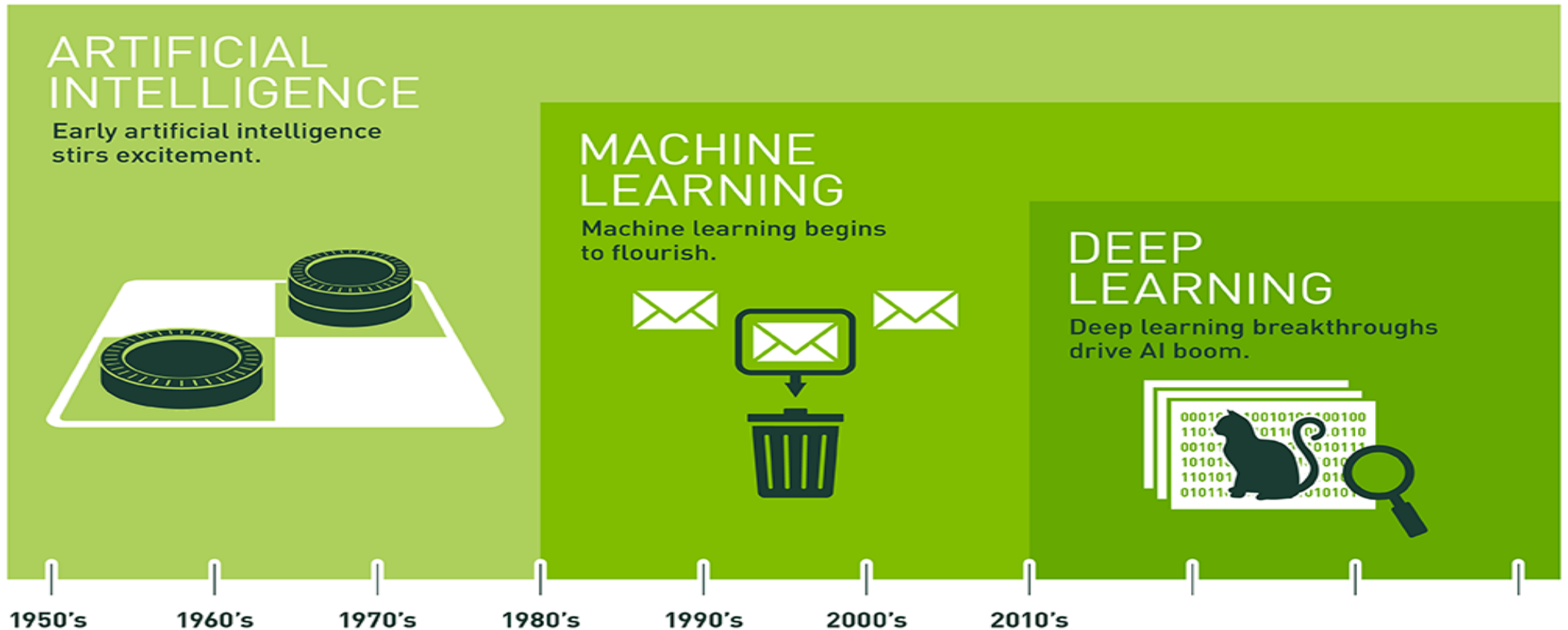e.g. Identify face

**Training data**
10-100M images

**Network architecture**
~10s-100s of layers
1B parameters

**Learning algorithm**
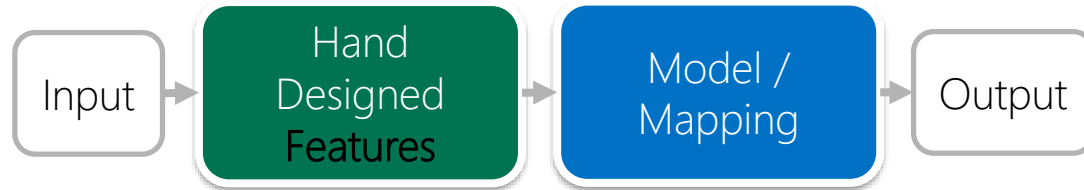~30 Exaflops
1-30 GPU days

# WHAT IS DEEP LEARNING?

# Accomplishing complex goals

# Difference in Workflow

Classic Machine Learning [ 1990 : now ]

Input → **Hand Designed Features** → **Model / Mapping** → Output

Examples [ Regression and SVMs ]



Deep/End-to-End Learning [ 2012 : now ]

Input → **Simple Features** → **Complex Features** → **Model/ Mapping** → Output

Example [ Conv Net ]

# Traditional Workflow

Classic Machine Learning [ 1990 : now ]

Input → Hand Designed **Features** → Model / Mapping → Output

Examples [ Regression and SVMs ]



**Challenge in Slack channel: How would you describe this image to someone (or something) blind?**

Difficult: From it's raw pixels.
Medium: From geometric primitives (lines, curves, colors)
Easy: Using any words that you may know

# Deep Learning Workflow



Experience: <u>Trust</u> Neural Network to learn features and model by providing inputs and outputs.

Key Skill: Experience (data) creation

Deep/End-to-End Learning [ 2012 : now ]

Input → Simple Features → Complex Features → Model/Mapping → Output

Example [ Conv Net ]

# NVIDIA'S DIGITS

# NVIDIA'S DIGITS
## Interactive Deep Learning GPU Training System

- Simplifies common deep learning tasks such as:

  - Managing data

  - Designing and training neural networks on multi-GPU systems

  - Monitoring performance in real time with advanced visualizations

- Completely interactive so data scientists can focus on designing and training networks rather than programming and debugging

- Open source

# NVIDIA'S DIGITS

## Interactive Deep Learning GPU Training System

Process Data | Configure DNN | Monitor Progress | Visualization

# DIGITS - MODEL

## New Object Detection Model

**Select Dataset** ❓

**Solver Options**

Training epochs ❓
30

Snapshot interval (in epochs) ❓
1

Validation interval (in epochs) ❓
1

Random seed ❓
[none]

Batch size ❓   multiples allowed
[network defaults]

Batch Accumulation ❓

Solver type ❓
Stochastic gradient descent (SGD)

Base Learning Rate ❓   multiples allowed
0.01

☐ Show advanced learning rate options

**Data Transformations**

Subtract Mean ❓
Image

Crop Size ❓
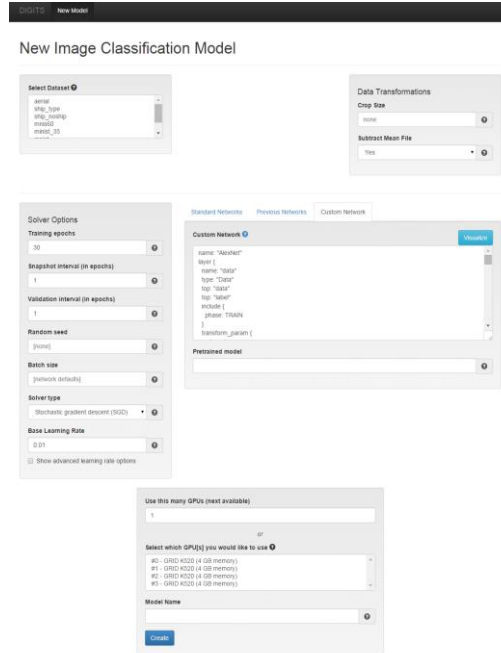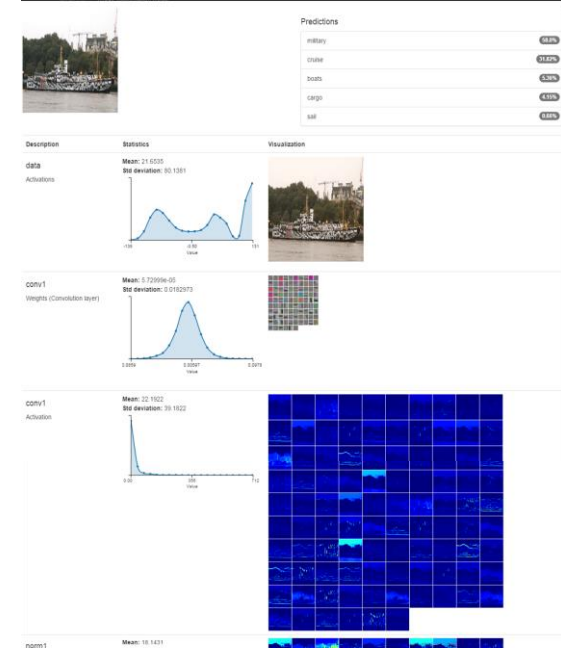none

**Python Layers** ❓
Server-side file ❓

☐ Use client-side file

Standard Networks | Previous Networks | Pretrained Networks | Custom Network

| Network | Details | Intended image size |
|---------|---------|---------------------|

## New Image Classification Model

**Select Dataset** ❓

**Solver Options**

Training epochs ❓
30

Snapshot interval (in epochs) ❓
1

Validation interval (in epochs) ❓
1

Random seed ❓
[none]

Batch size ❓   multiples allowed
[network defaults]

Batch Accumulation ❓

Solver type ❓
Stochastic gradient descent (SGD)

Base Learning Rate ❓   multiples allowed
0.01

☐ Show advanced learning rate options

**Data Transformations**

Subtract Mean ❓
Image

Crop Size ❓
none

**Python Layers** ❓
Server-side file ❓

☐ Use client-side file

Standard Networks | Previous Networks | Pretrained Networks | Custom Network

Caffe | Torch

| Network | Details | Intended image size |
|---------|---------|---------------------|
| ○ LeNet | Original paper [1998] | 28x28 (gray) |

Define custom layers with Python

Can anneal the learning rate

Differences may exist between model tasks

# DIGITS – VISUALIZATION RESULTS

# ENHANCING IMAGES WITH AN AI AUTOENCODER

# A great candidate for Deep Learning!

INPUT x

FUNCTION f:

OUTPUT f(x)

# Training Set of images.



1 sample per pixel

- It requires pairs of noisy and noise-free images. The network will learn to remove the noise from the images.

- We can then deploy this trained model to any image we want to denoise. (inference)

# Deep Learning for Image Denoising

**Training Data**

**Training**

**Trained Neural Network**

**Inferencing**

| Collect iamhes Add Noise to training images | Training on progression of images | Trained network detects noise and reconstructs |
|---|---|---|

Apply trained network to noisy images

# Learning about images (CNN)

Raw data Low-level features Mid-level features High-level features



Input

Result

**Application components:**

**Task objective**
e.g. Identify face

**Training data**
10-100M images

**Network architecture**
~10s-100s of layers
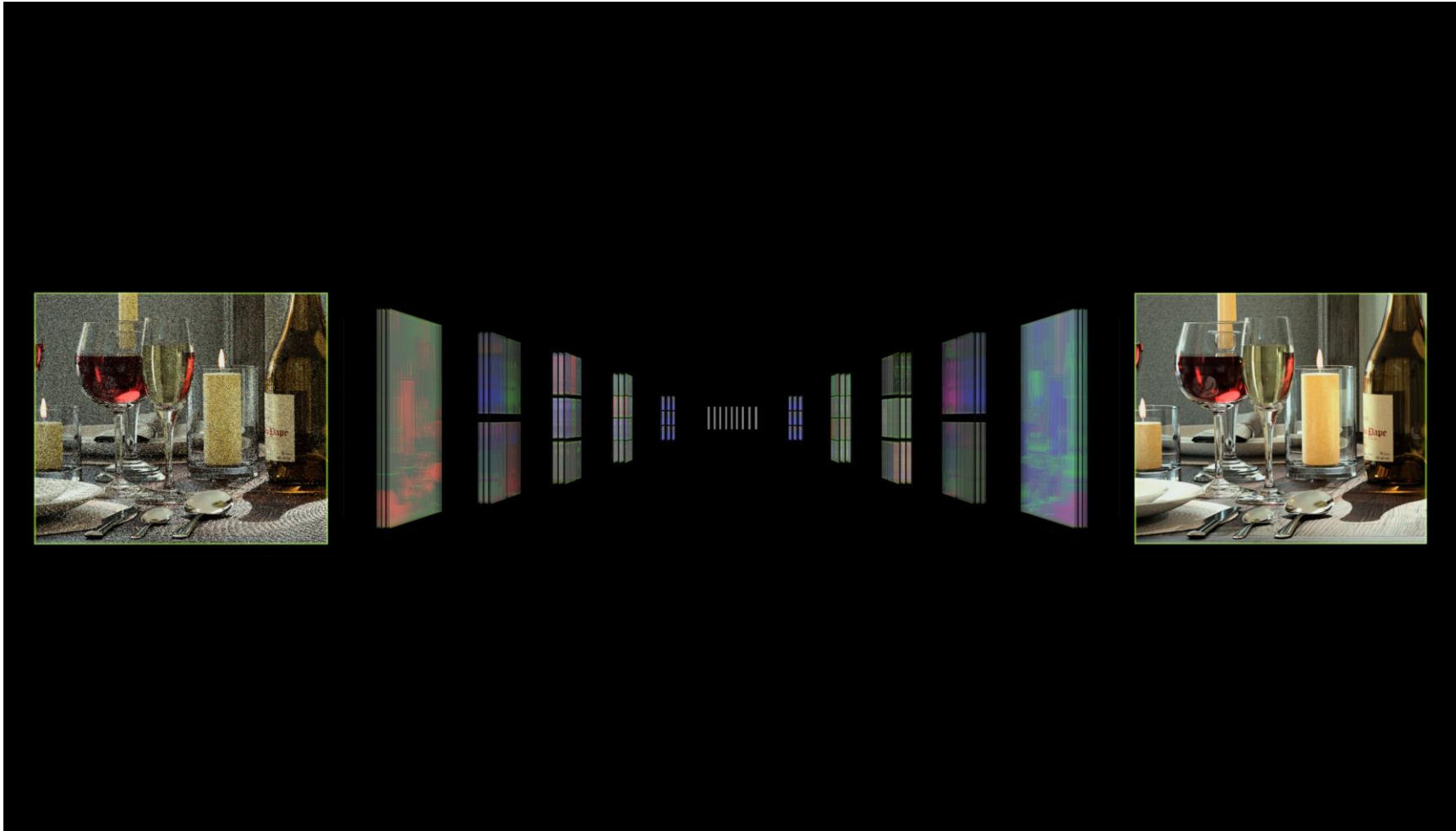1B parameters

**Learning algorithm**
~30 Exaflops
1-30 GPU days

# Autoencoder – in Action

Apply Noise
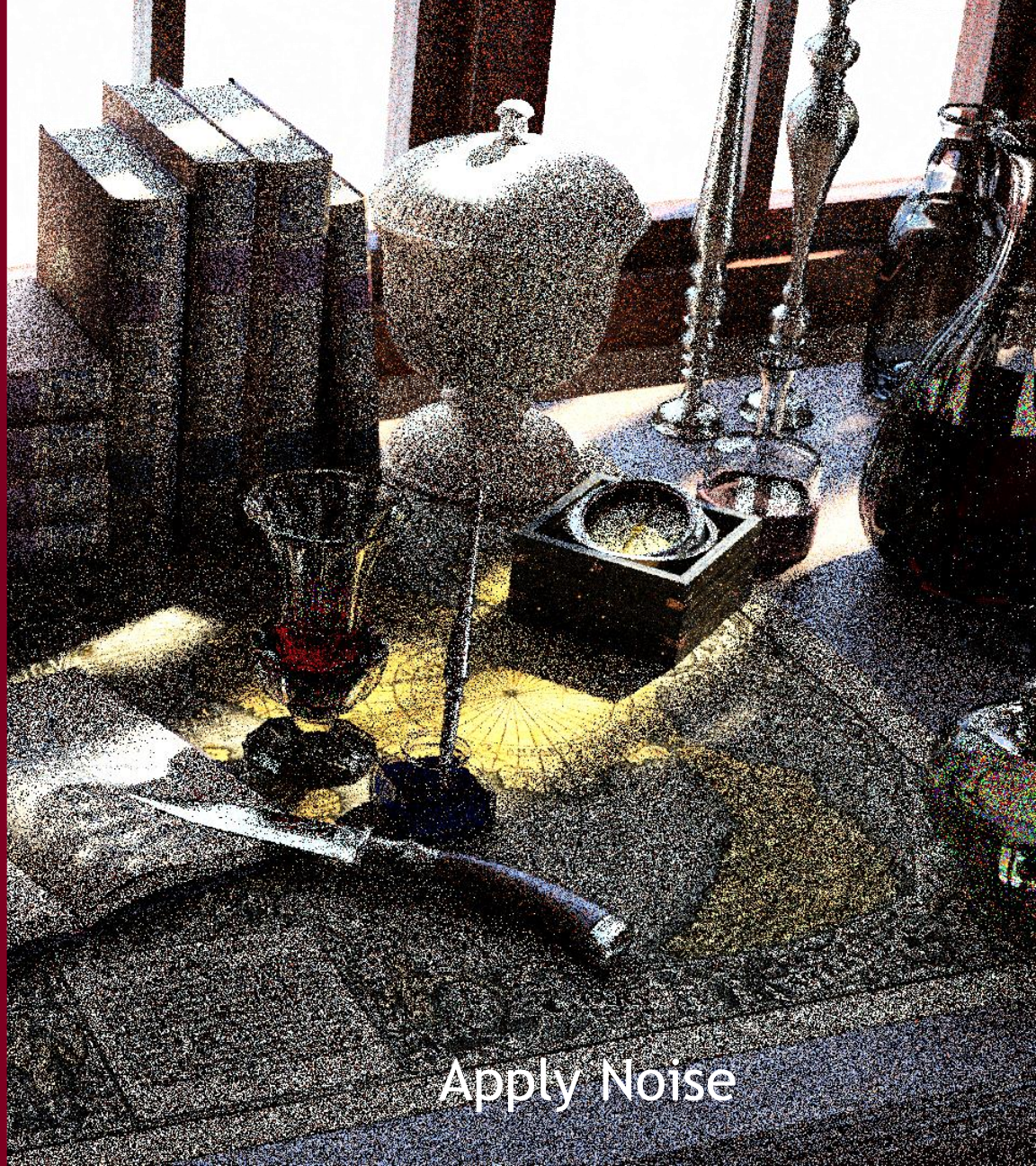Apply Noise

Provide image to autoencoder enhance

# TensorRT

# SOFTWARE INFERENCING PERFROMANCE EHANCEMENT

# NVIDIA DEEP LEARNING SOFTWARE PLATFORM



**TRAINING FRAMEWORK**

Training Data

Data Management

Training

Model Assessment

Trained Neural Network

**TensorRT**

Embedded

Automotive

Data center

**NVIDIA DEEP LEARNING SDK**

# NVIDIA TensorRT

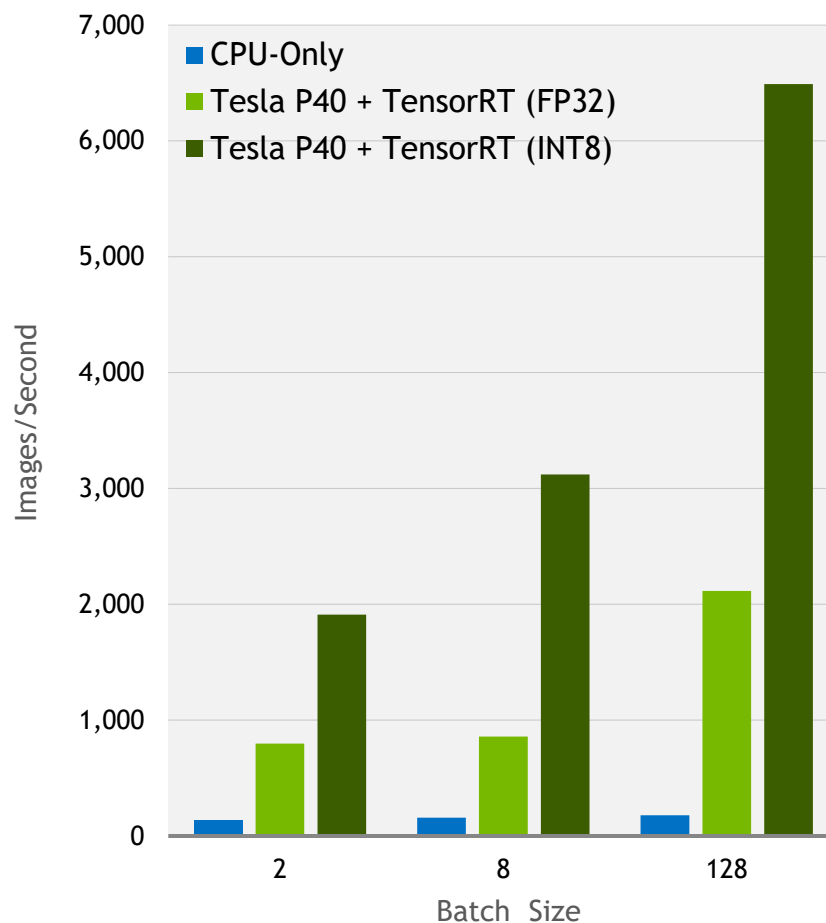High-performance deep learning inference for production deployment

High performance neural network inference engine for production deployment

Generate optimized and deployment-ready models for datacenter, embedded and automotive platforms

Deliver high-performance, low-latency inference demanded by real-time services

Deploy faster, more responsive and memory efficient deep learning applications with INT8 and FP16 optimized precision support

developer.nvidia.com/tensorrt

## Up to 36x More Image/sec

Legend:
- CPU-Only
- Tesla P40 + TensorRT (FP32)
- Tesla P40 + TensorRT (INT8)

Y-axis: Images/Second (0 to 7,000)
X-axis: Batch Size (2, 8, 128)

*GoogLenet, CPU-only vs Tesla P40 + TensorRT*
*CPU: 1 socket E4 2690 v4 @2.6 GHz, HT-on*
*GPU: 2 socket E5-2698 v3 @2.3 GHz, HT off, 1 P40 card in the box*

31

# TENSORRT

## Networks Supported

- **Image Classification (AlexNet, GoogleNet, VGG, ResNet)**
- **Object Detection**
- **Segmentation**

## Not Yet Supported

- **RNN/LSTM**
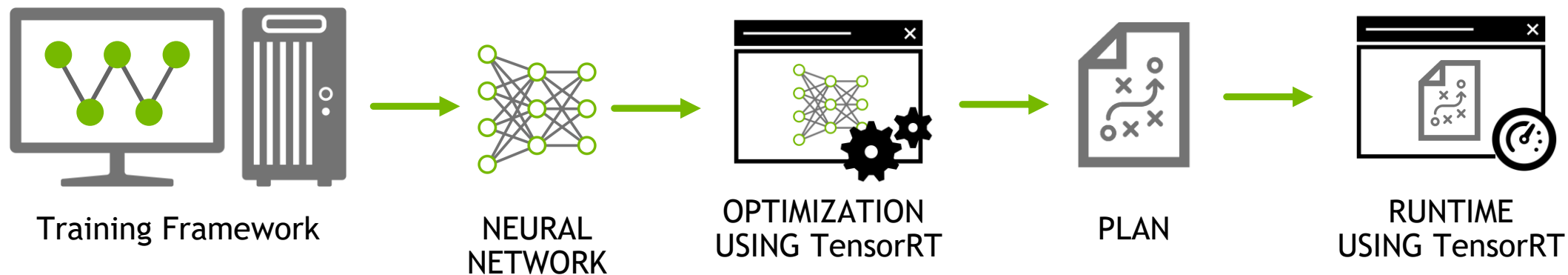- **3D convolutions**
- **Custom user layers**

# TENSORRT
## Layers Types Supported

- **Convolution:** Currently only 2D convolutions
- **Activation:** ReLU, tanh and sigmoid
- **Pooling:** max and average
- **Scale:** similar to Caffe Power layer (shift+scale*x)^p
- **ElementWise:** sum, product or max of two tensors
- **LRN:** cross-channel only
- **Fully-connected:** with or without bias
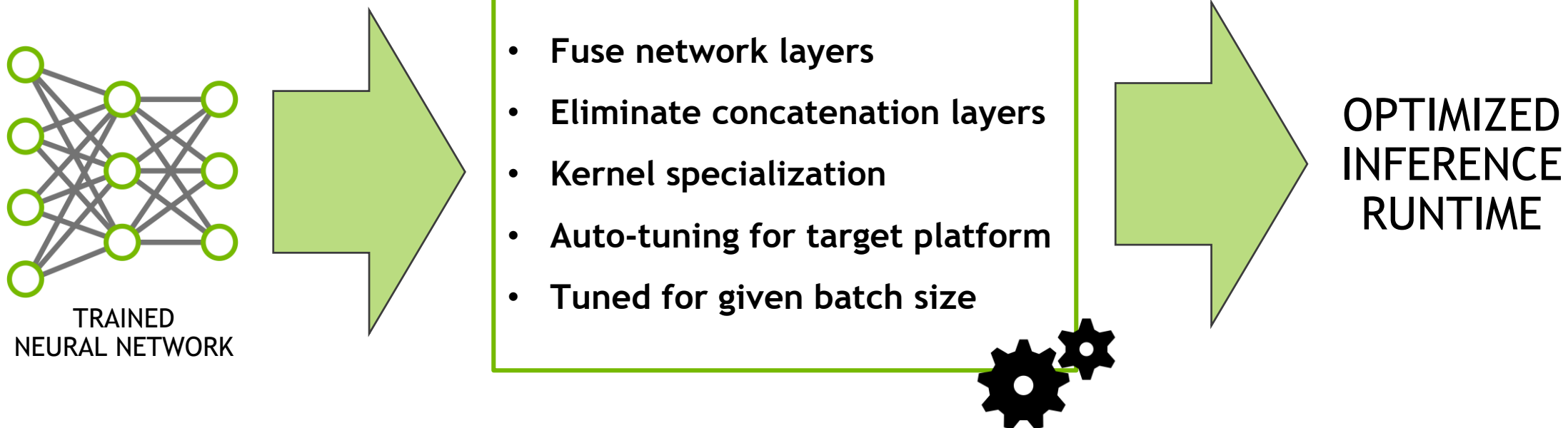- **SoftMax:** cross-channel only
- **Deconvolution**
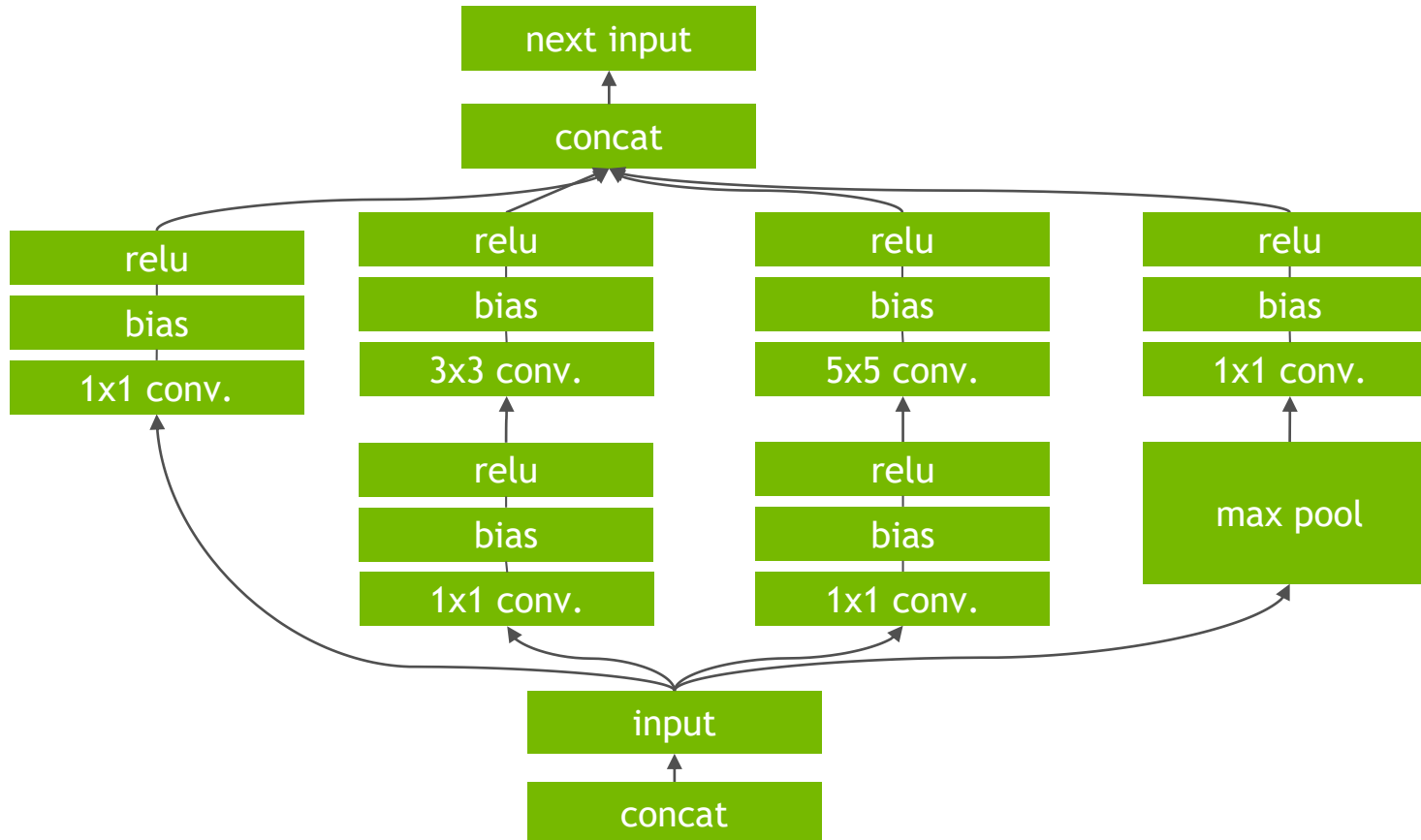
# TENSORRT
## Workflow

**Training Framework** → **NEURAL NETWORK** → **OPTIMIZATION USING TensorRT** → **PLAN** → **RUNTIME USING TensorRT**

# TENSORRT
## Optimizations

TRAINED
NEURAL NETWORK

- **Fuse network layers**
- **Eliminate concatenation layers**
- **Kernel specialization**
- **Auto-tuning for target platform**
- **Tuned for given batch size**

OPTIMIZED
INFERENCE
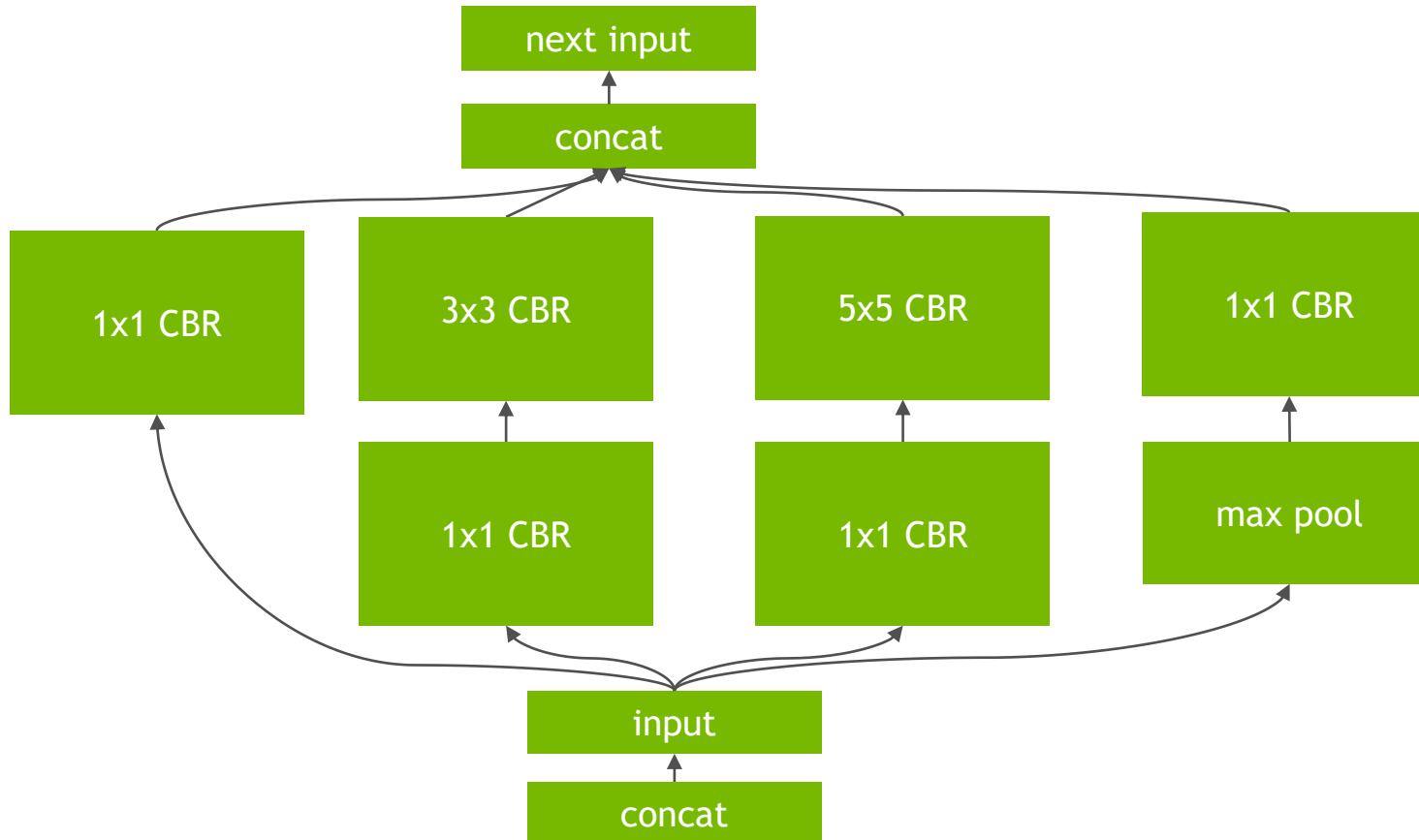RUNTIME

# GRAPH OPTIMIZATION
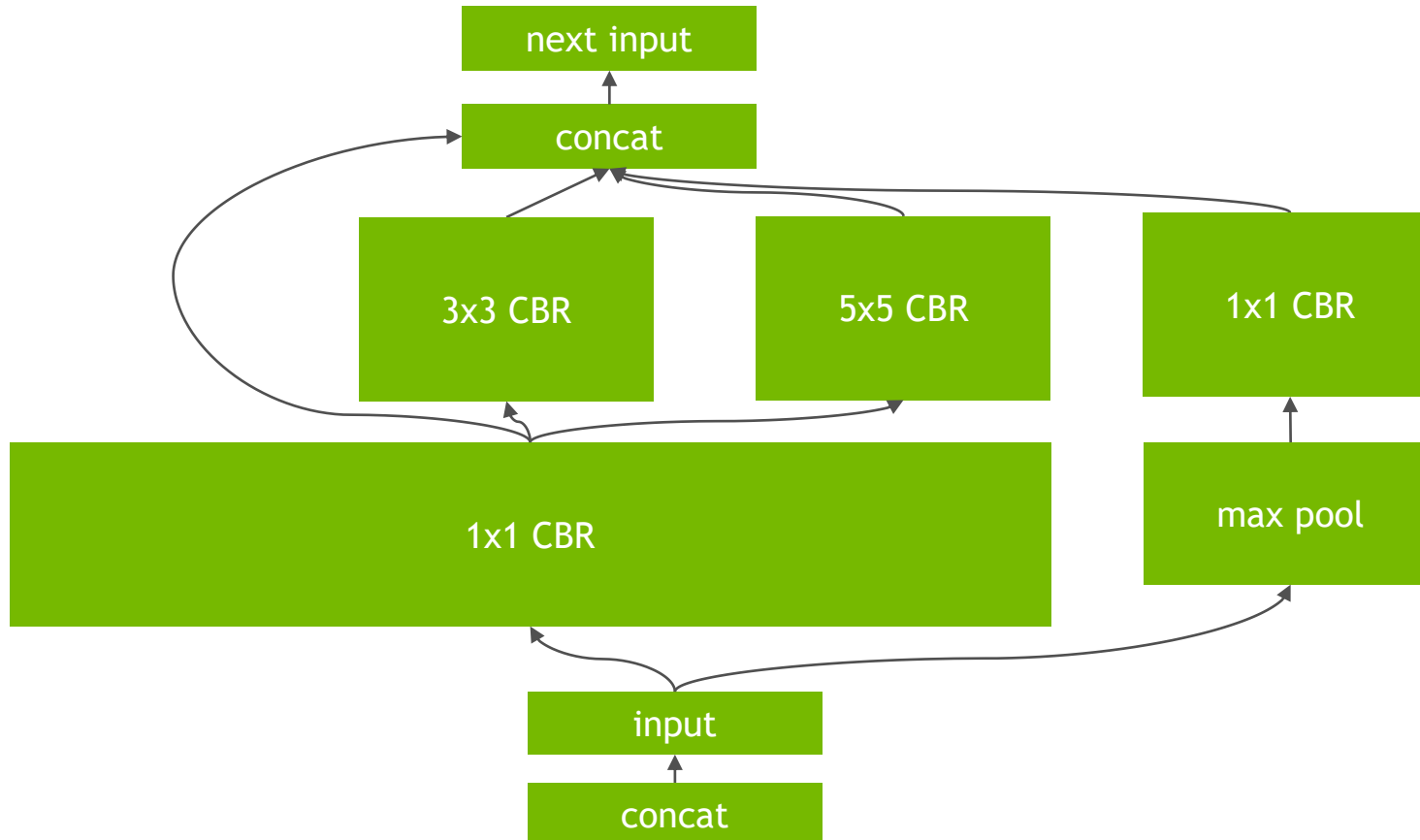## Unoptimized network

# GRAPH OPTIMIZATION

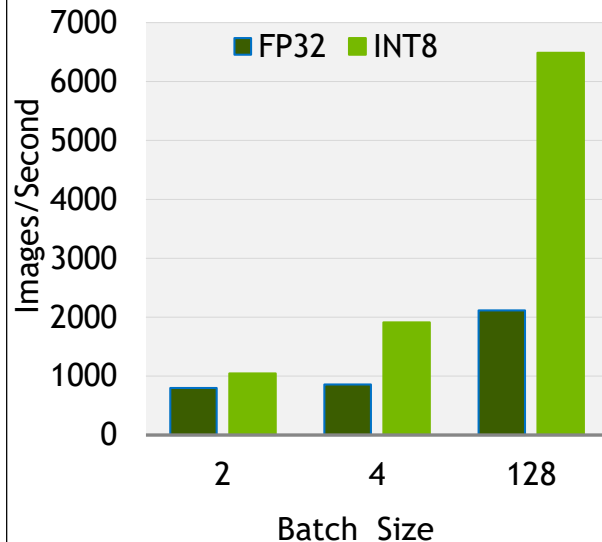## Vertical fusion

# GRAPH OPTIMIZATION
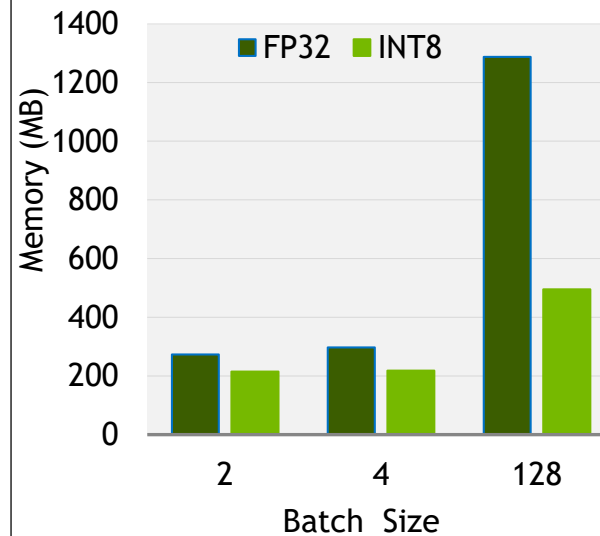
## Horizontal fusion

# INT8 PRECISION
## New in TensorRT

### PERFORMANCE
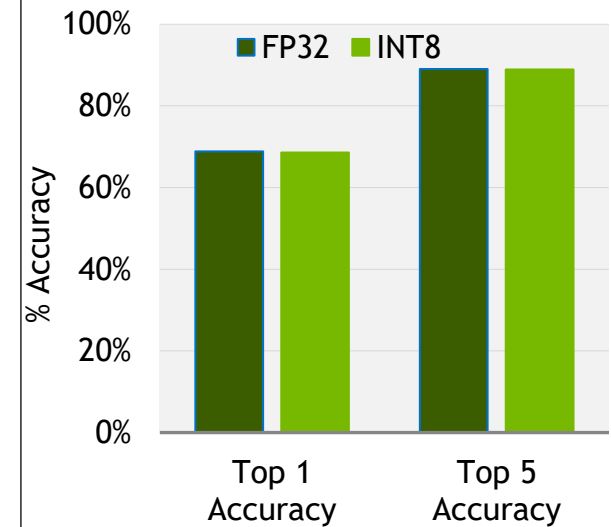
Up To 3x More Images/sec with INT8 Precision



Images/Second vs Batch Size (FP32, INT8)

### EFFICIENCY

Deploy 2x Larger Models with INT8 Precision



Memory (MB) vs Batch Size (FP32, INT8)

### ACCURACY

Deliver full accuracy with INT8 precision



% Accuracy: Top 1 Accuracy, Top 5 Accuracy (FP32, INT8)

*GoogLenet*, *FP32 vs INT8 precision + TensorRT on*
*Tesla P40 GPU*, *2 Socket Haswell E5-2698 v3 @2.3GHz with HT off*

# THANK YOU