

# HIGH-PERFORMANCE INPUT PIPELINES FOR SCALABLE DEEP LEARNING

**Miroslav Klivansky**  
**Pure Storage**



**school of ai**



# **QUESTION ON EVERYONE'S MIND:**

# **WHY IS A STORAGE COMPANY HERE?**



“We don’t have better algorithms,  
we just have more data”

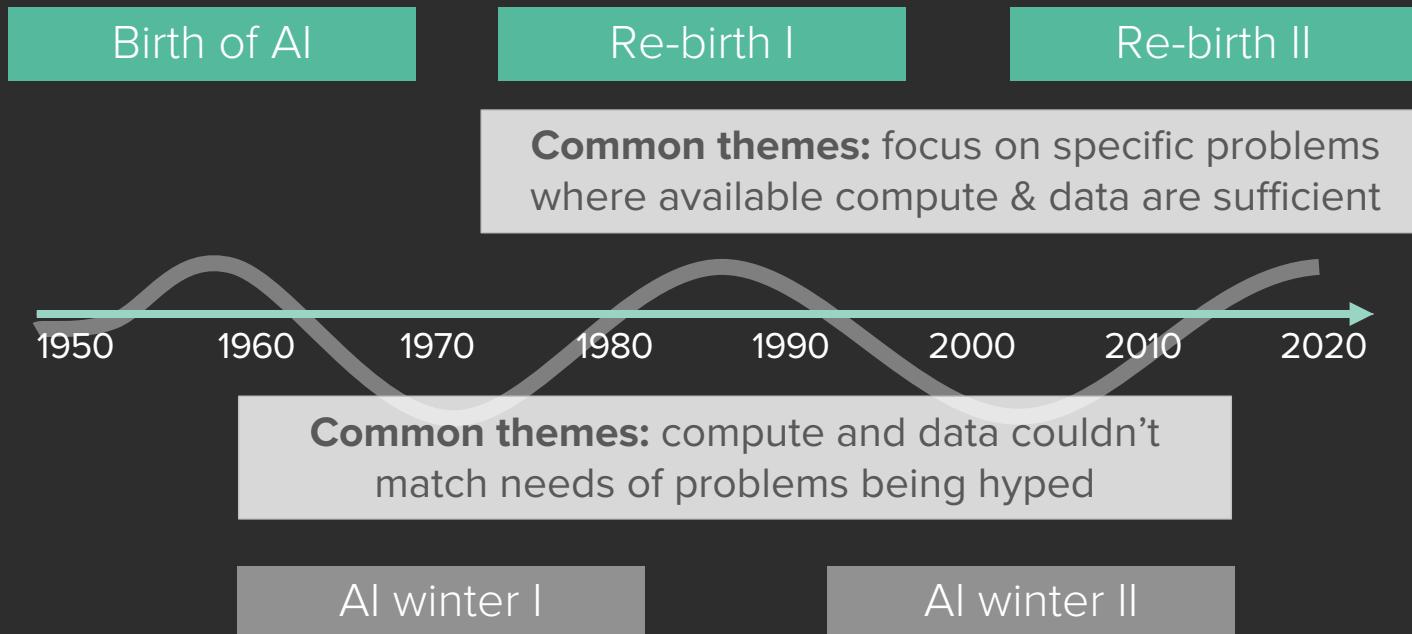
**PETER NORVIG**  
Engineering Director, Google

# The AI “Hierarchy of Needs”

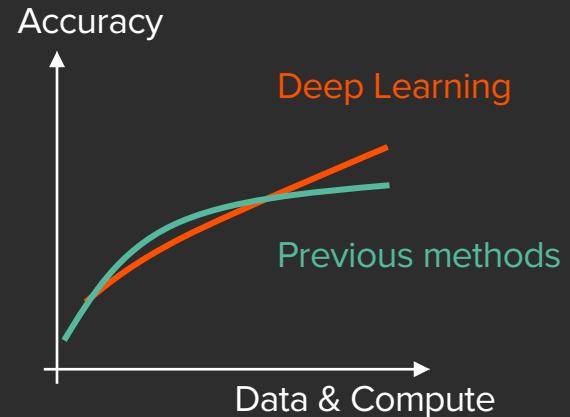
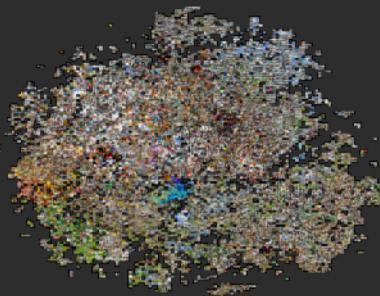
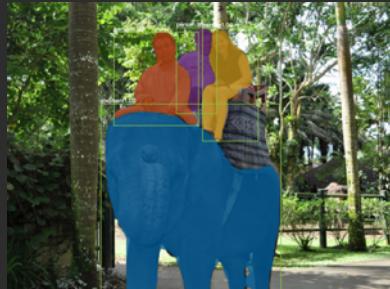
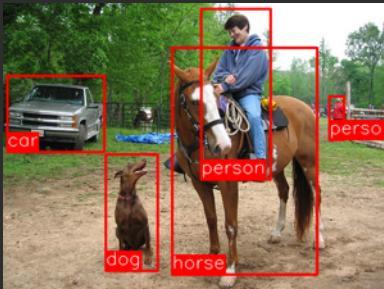
*credit: Monica Rogati*



# THIS IS NOT THE FIRST AI HYPE WAVE



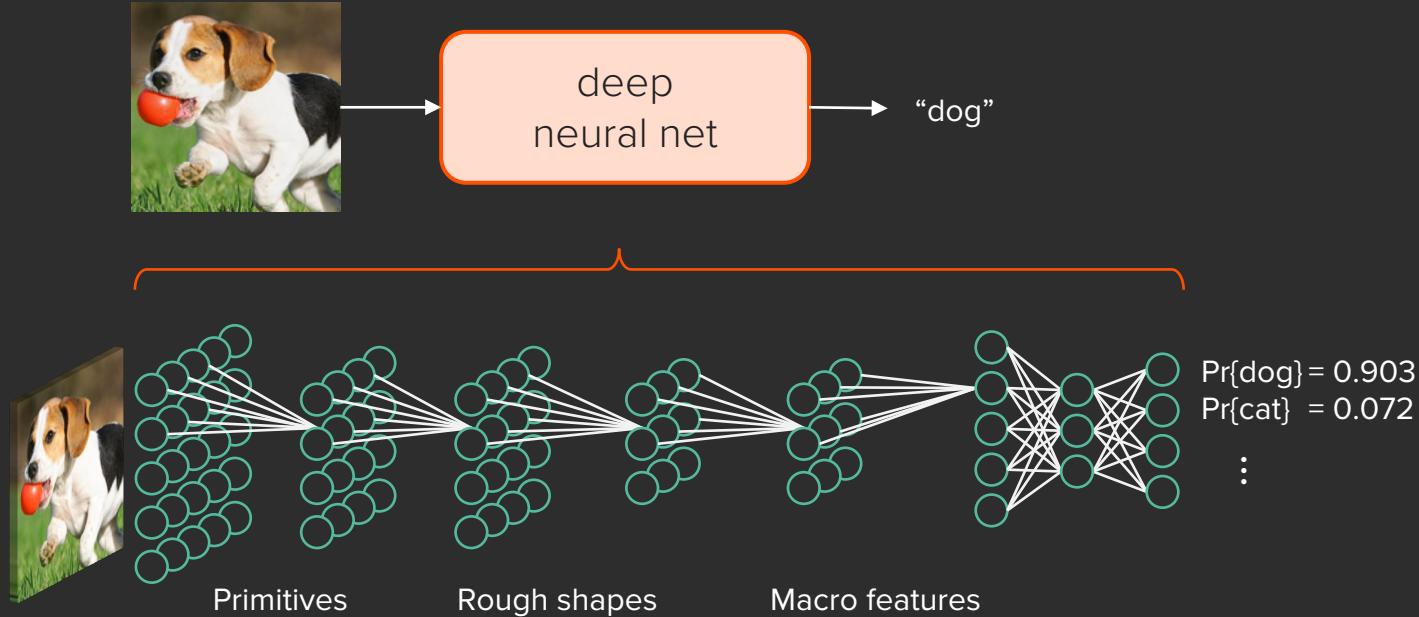
# DEEP LEARNING = MASSIVE DATA & COMPUTE



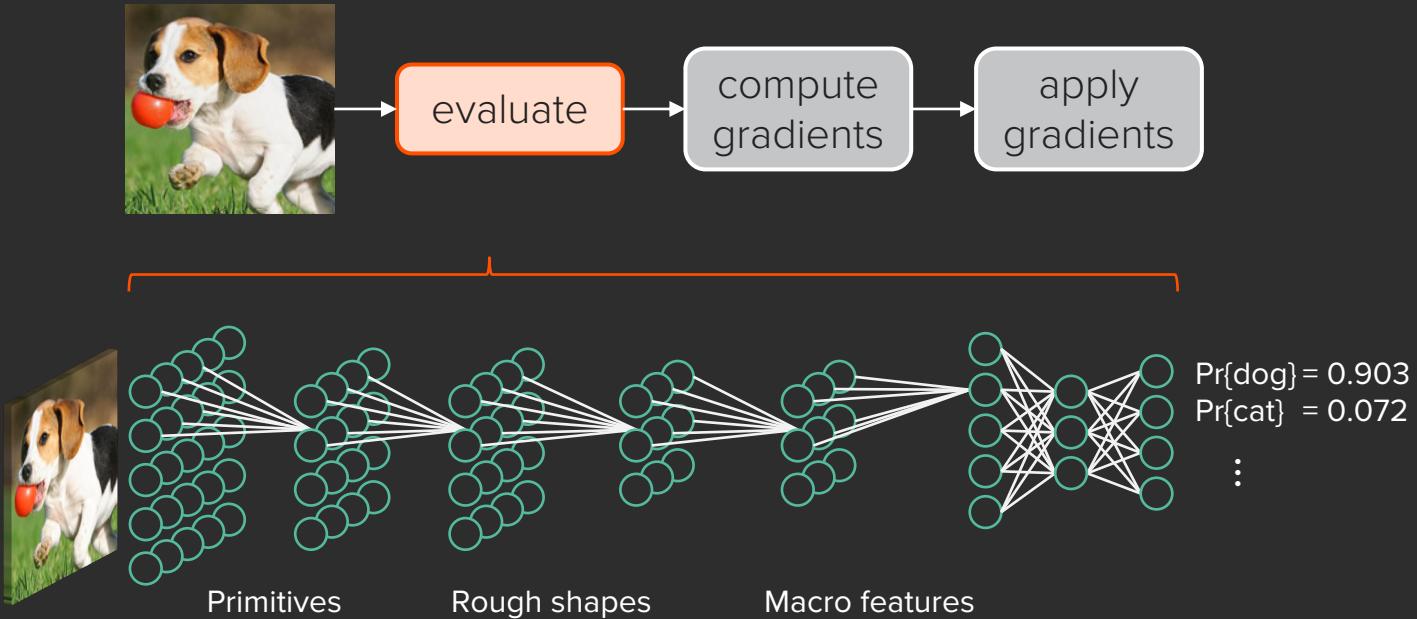
**STATE-OF-THE-ART RESULTS ACROSS VISION, SPEECH, LANGUAGE, AND MORE**

# THE INTUITION BEHIND DEEP LEARNING

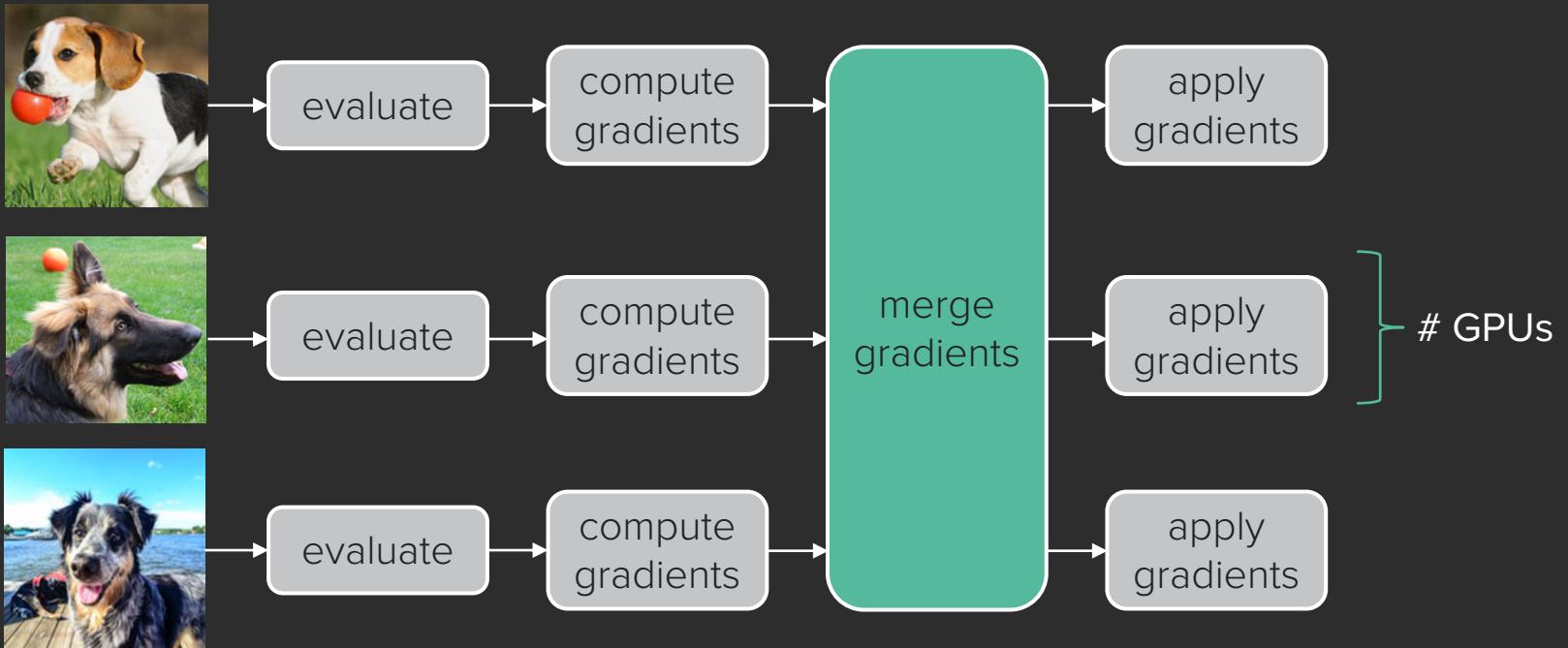
---



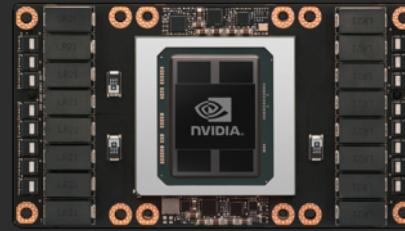
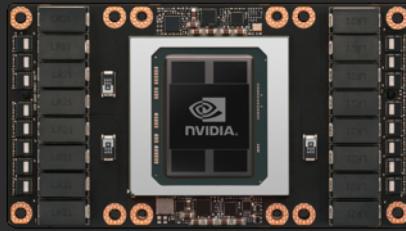
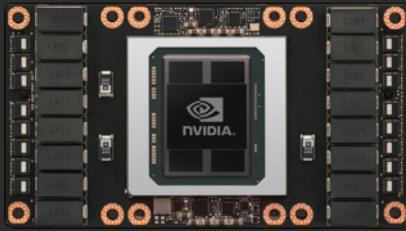
# TRAINING A DEEP NEURAL NETWORK



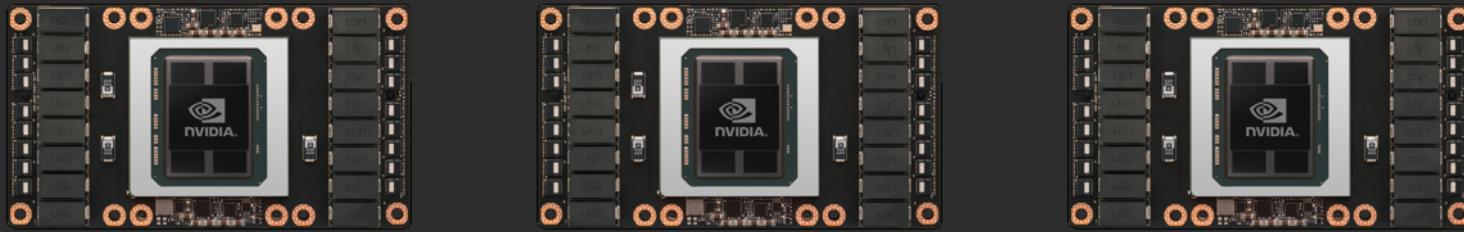
# DISTRIBUTED TRAINING



# MORE, FASTER GPUs + MORE DATA



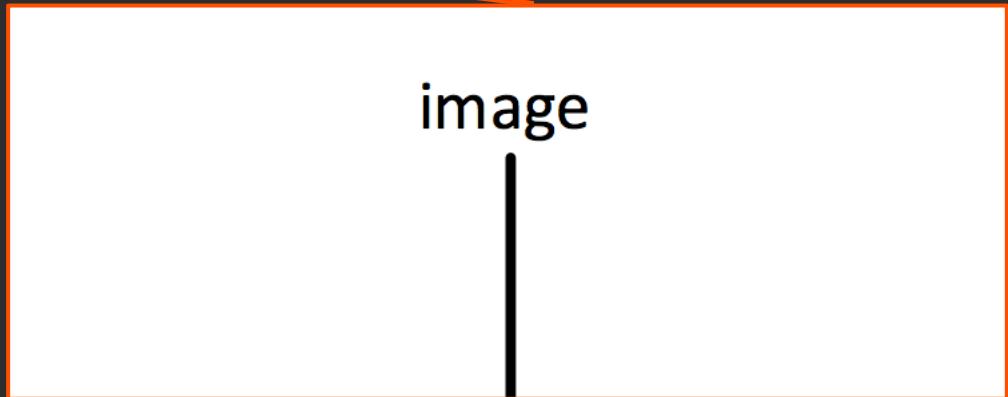
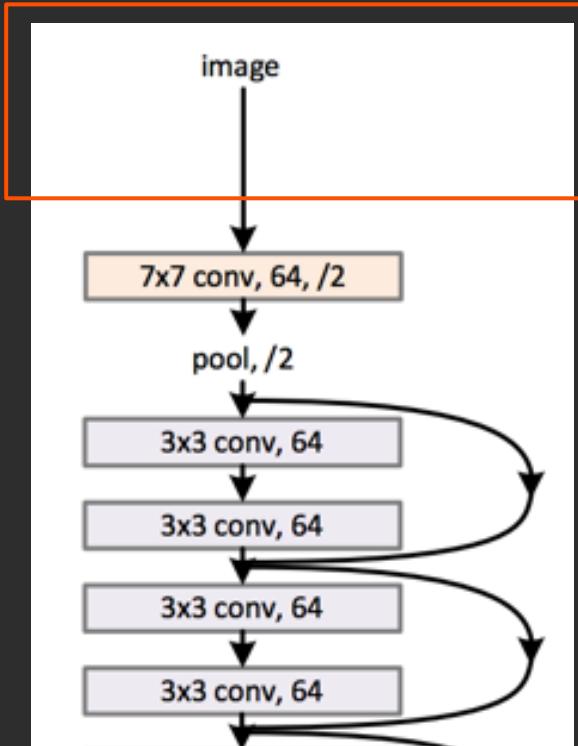
# CAN WE KEEP GPUs FED WITH DATA?



**INPUT PIPELINE = POTENTIAL BOTTLENECK**

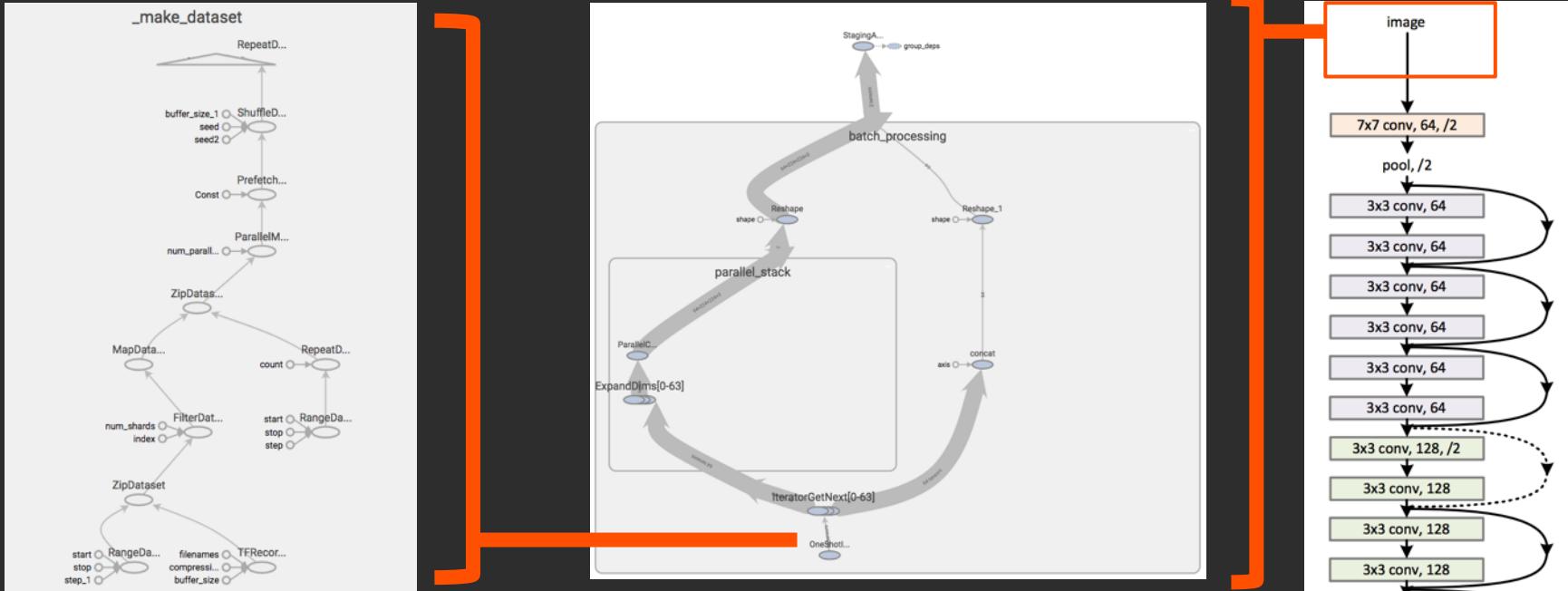


# INPUT PIPELINES



CAN IT BE THAT SIMPLE?

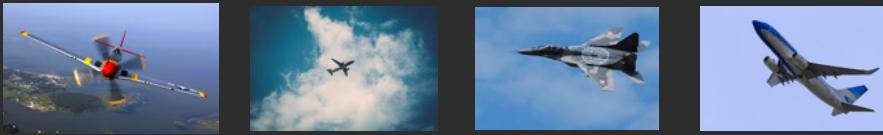
# REAL INPUT PIPELINES



# CAN YOU SPOT THE BOTTLENECK?

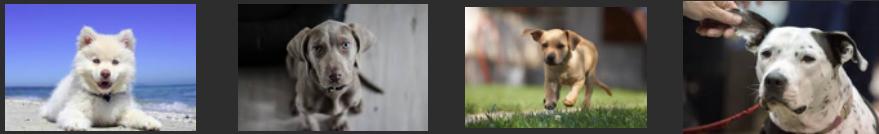
# FROM IMAGES TO TENSORS

PLANE



1. Enumerate

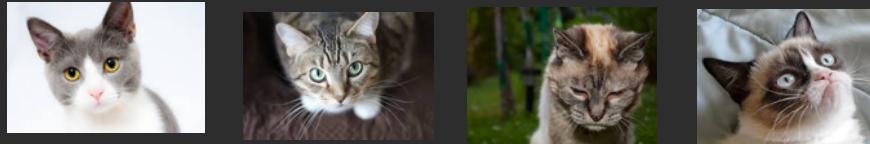
DOG



BOAT



CAT



# FROM IMAGES TO TENSORS



PLANE



PLANE



PLANE



PLANE



DOG



DOG



DOG



DOG



BOAT



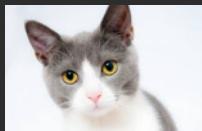
BOAT



BOAT



BOAT



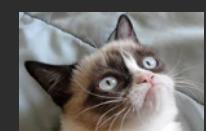
CAT



CAT



CAT



CAT

1. Enumerate
2. Associate labels

# FROM IMAGES TO TENSORS



BOAT



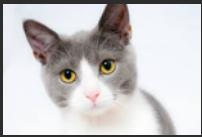
CAT



BOAT



BOAT



CAT



PLANE



PLANE



CAT



PLANE



DOG



DOG



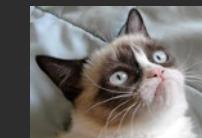
DOG



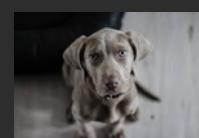
BOAT



PLANE



CAT



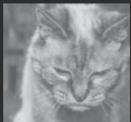
DOG

1. Enumerate
2. Associate labels
3. Shuffle

# FROM IMAGES TO TENSORS



BOAT



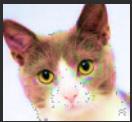
CAT



BOAT



BOAT



CAT



PLANE



PLANE



CAT



PLANE



DOG



DOG



BOAT



PLANE



CAT



DOG

1. Enumerate
2. Associate labels
3. Shuffle
4. Read, crop, distort

# FROM IMAGES TO TENSORS



BOAT



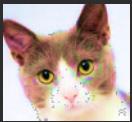
CAT



BOAT



BOAT



CAT



PLANE



PLANE



CAT



PLANE



DOG



DOG



DOG



BOAT



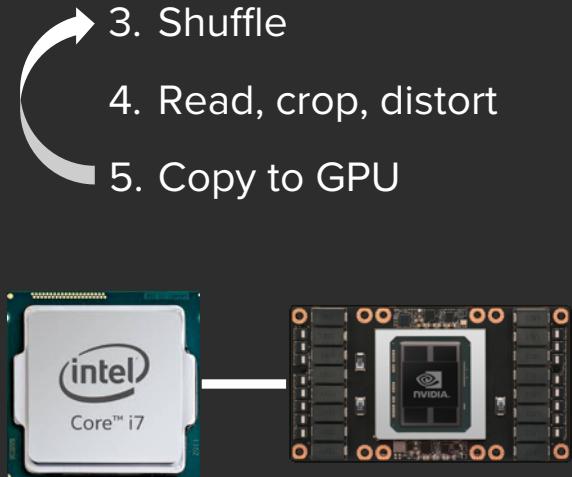
PLANE



CAT



DOG



# FROM IMAGES TO TENSORS

Other domains (NLP, speech, etc.)  
will follow a similar(ish) flow

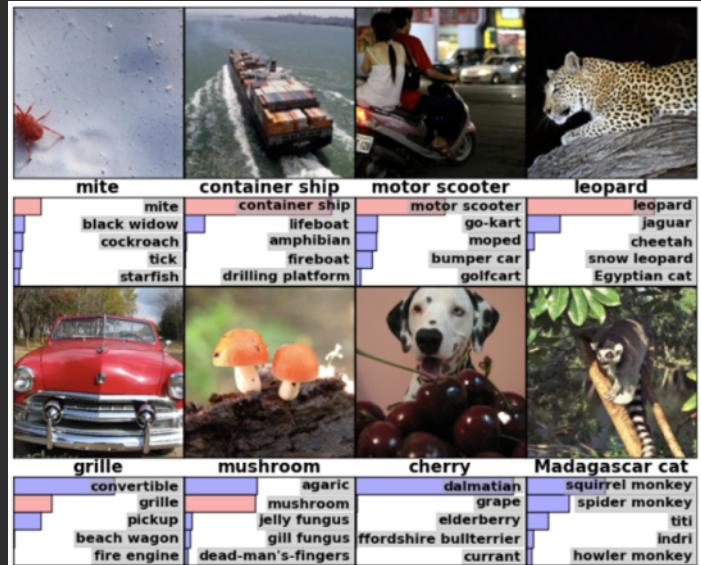
**ANY OF THESE STEPS CAN BE  
A POTENTIAL BOTTLENECK**

1. Enumerate
2. Associate labels
3. Shuffle
4. Read, crop, distort
5. Copy to GPU

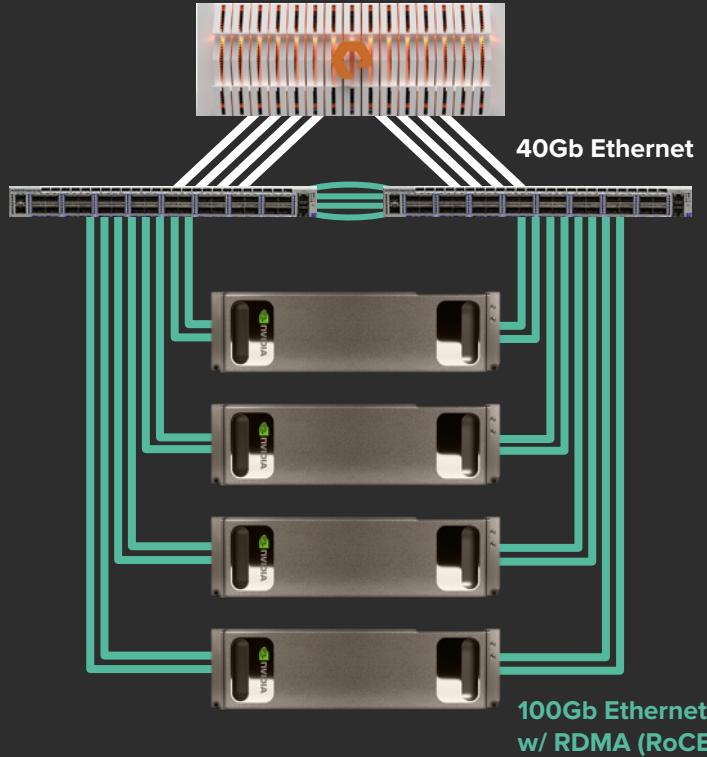
# EVALUATION METHODOLOGY



1.3M images, 1000 categories



# HARDWARE STACK

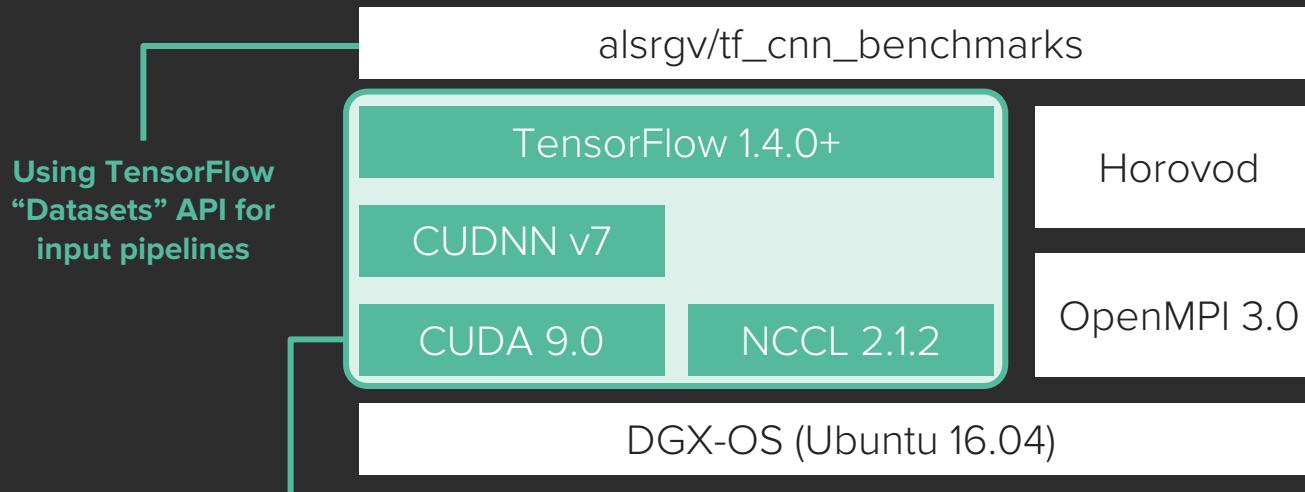


**Pure Storage FlashBlade: 15x17TB**  
179T usable before data reduction

**Arista DCS-7060CX2-32S**  
32x 100Gb/s QSFP100 ports

**4x NVIDIA DGX-1, each with**  
8x Tesla V100 GPUs (SXM2)  
2x Intel E5-2698 v4 @ 2.20GHz  
4x Mellanox MT27700 100Gb/s VPI adapters  
512GB DDR4-2400

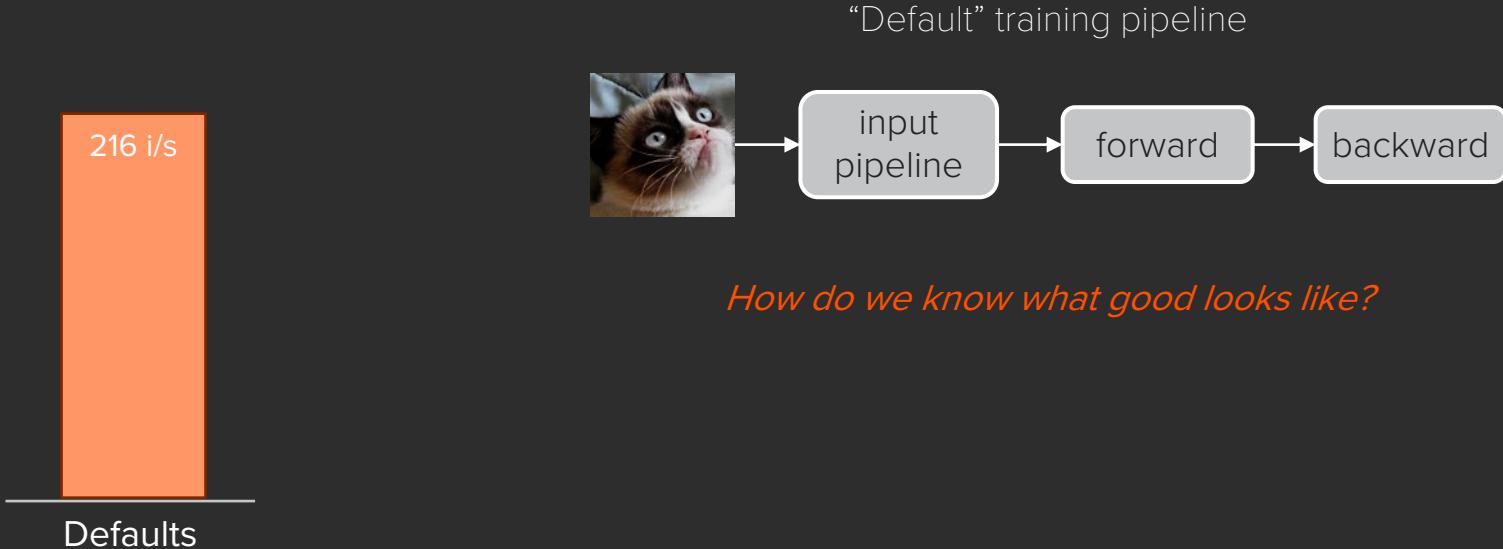
# SOFTWARE STACK



[nvcr.io/nvidia/tensorflow:17.12](https://nvcr.io/nvidia/tensorflow:17.12)

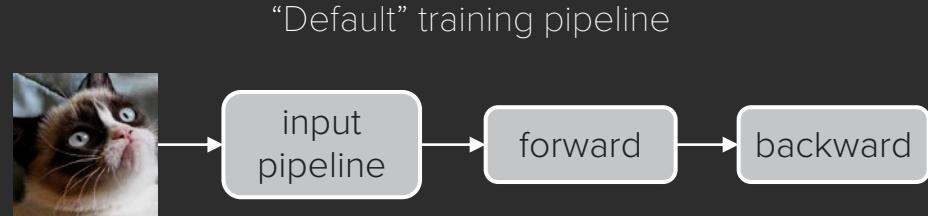
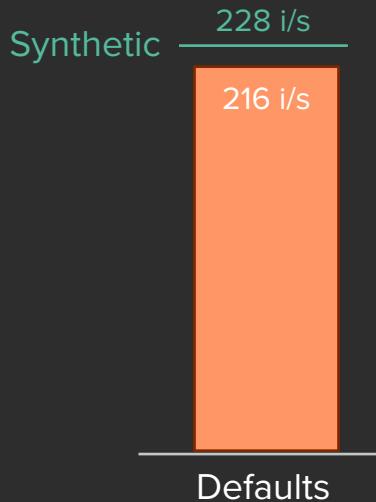
# TRAINING WITH 1 GPU

*Images per second when training Inception3 (batch size = 64)*



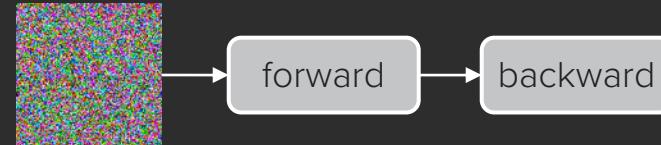
# TRAINING WITH 1 GPU

*Images per second when training Inception3 (batch size = 64)*



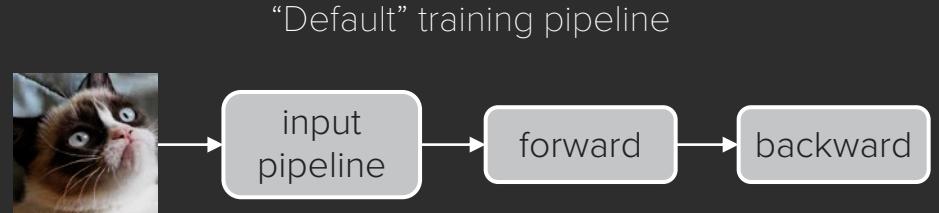
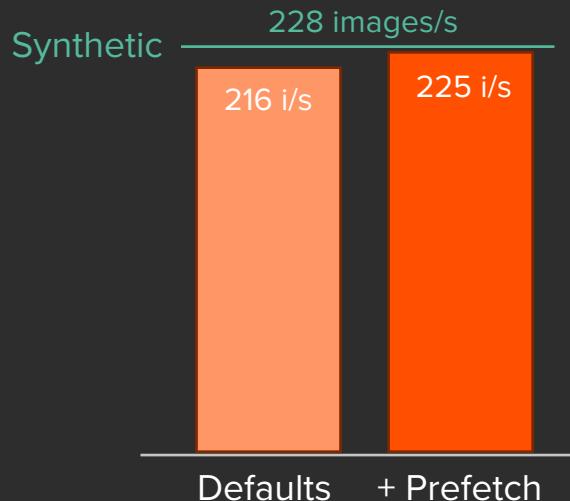
*How do we know what good looks like?*

Replace the input pipeline with synthetic data

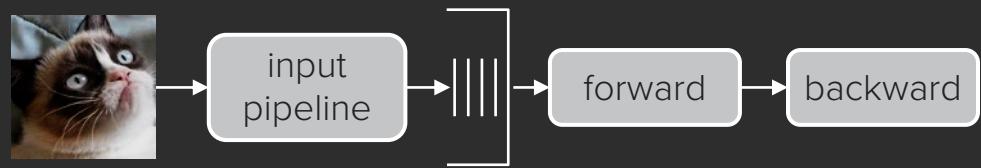


# TRAINING WITH 1 GPU

*Images per second when training Inception3 (batch size = 64)*



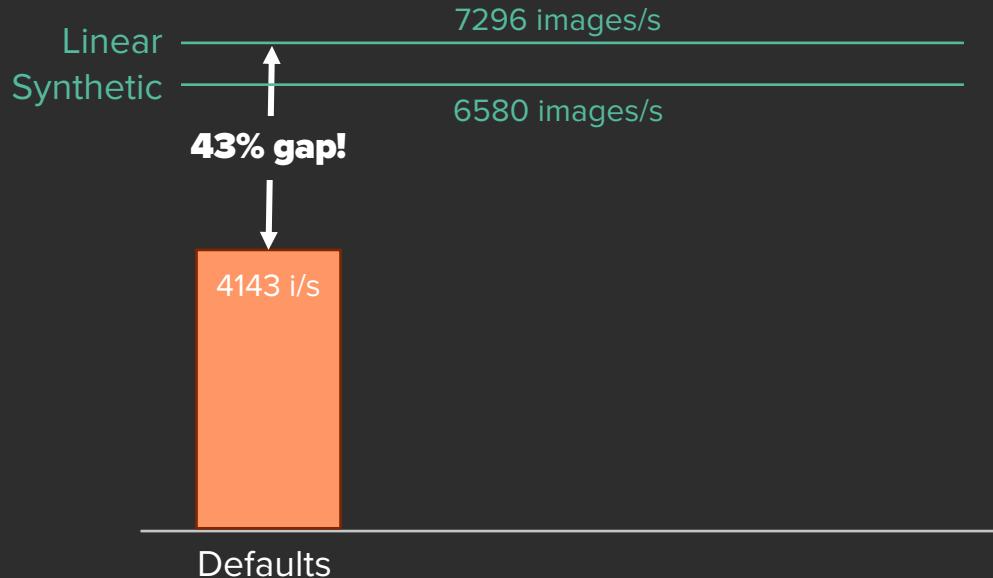
Adding a prefetch queue improves scheduler behavior



**SHOULD WE CARE ABOUT ~1.3%?**

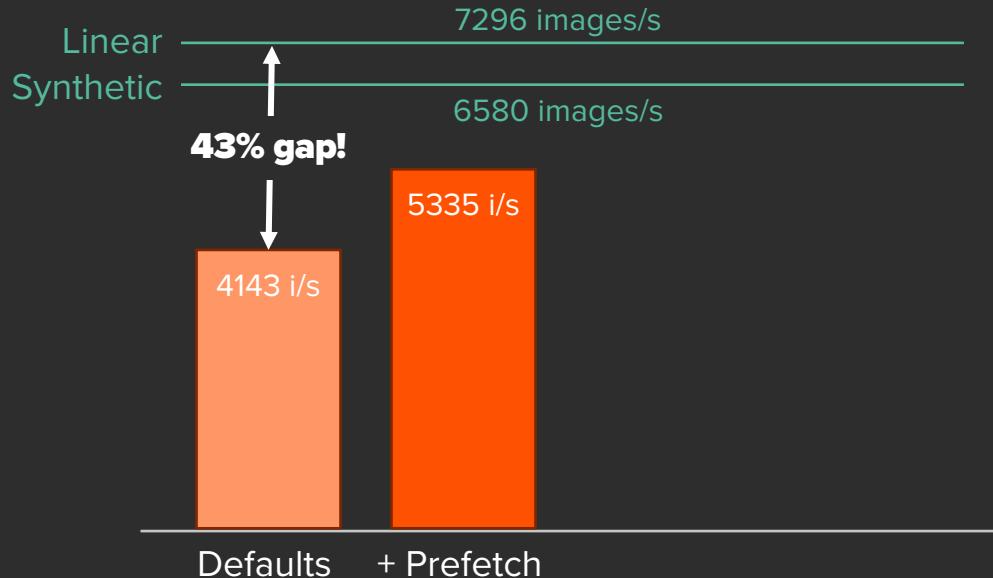
# SCALING TO 32 GPUs (4x DGX-1s)

*Images per second when training Inception3 (batch size = 64/GPU)*



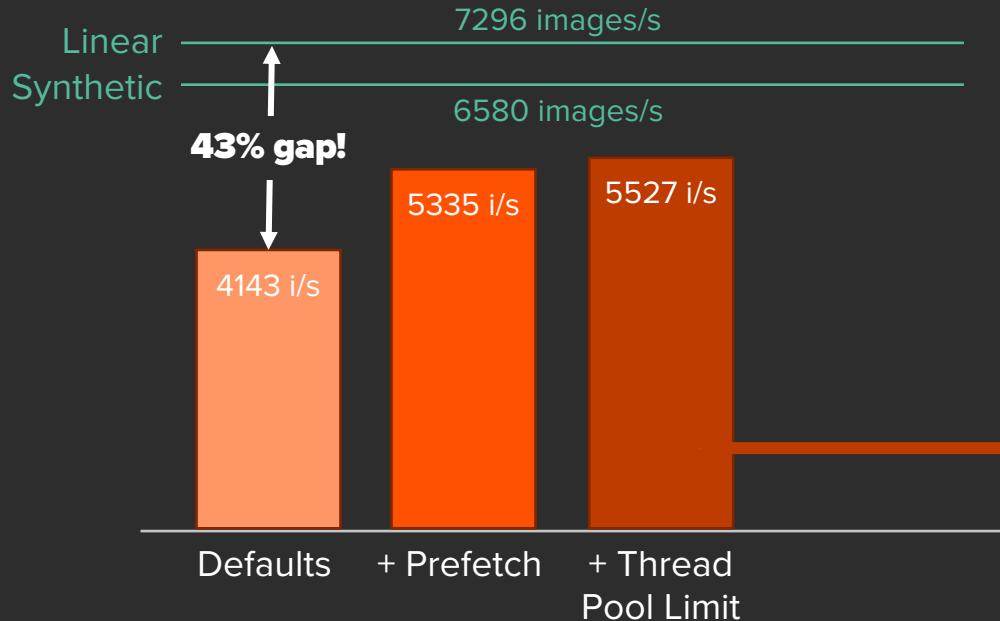
# SCALING TO 32 GPUs (4x DGX-1s)

*Images per second when training Inception3 (batch size = 64/GPU)*



# SCALING TO 32 GPUs (4x DGX-1s)

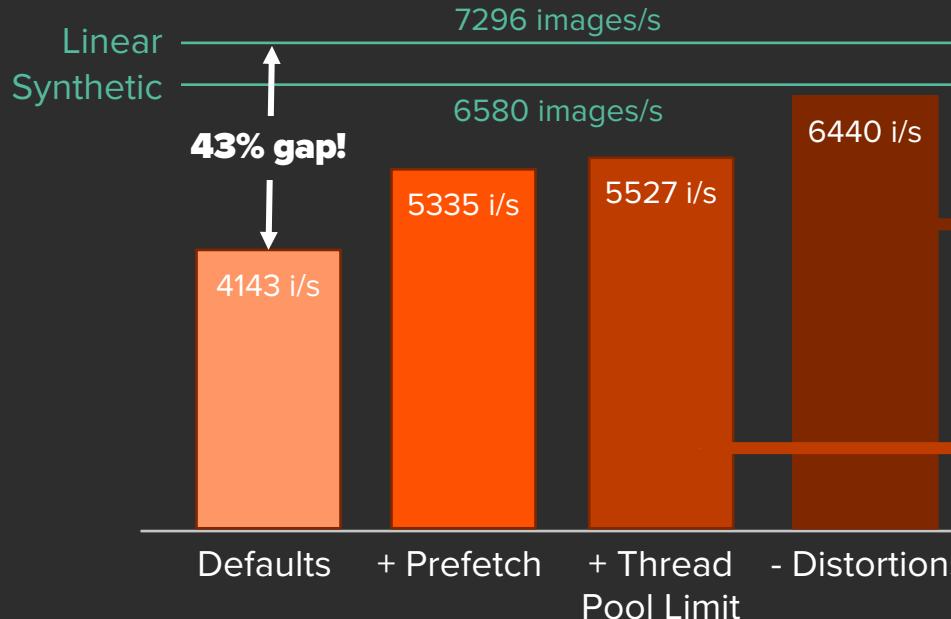
*Images per second when training Inception3 (batch size = 64/GPU)*



**Thread pool limits:** Avoid over-subscribing CPU with too many threads.  
(`inter_op_parallelism_threads`)

# SCALING TO 32 GPUs (4x DGX-1s)

*Images per second when training Inception3 (batch size = 64/GPU)*



**No Distortions:** Skip preprocessing step from input pipeline. This is an unrealistic configuration, but it shows the bottleneck.

**Thread pool limits:** Avoid over-subscribing CPU with too many threads.  
(`inter_op_parallelism_threads`)

# 2.5X Performance Improvement

