

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

LAB 1, PART 1: INTRODUCTION AND MOTIVATION



DEEP  
LEARNING  
INSTITUTE

# COURSE OVERVIEW

- Lab 1: Gradient Descent vs Stochastic Gradient Descent, and the Effects of Batch Size
- Lab 2: Multi-GPU DL Training Implementation using Horovod
- Lab 3: Algorithmic Concerns for Training at Scale

# LAB 1 OVERVIEW

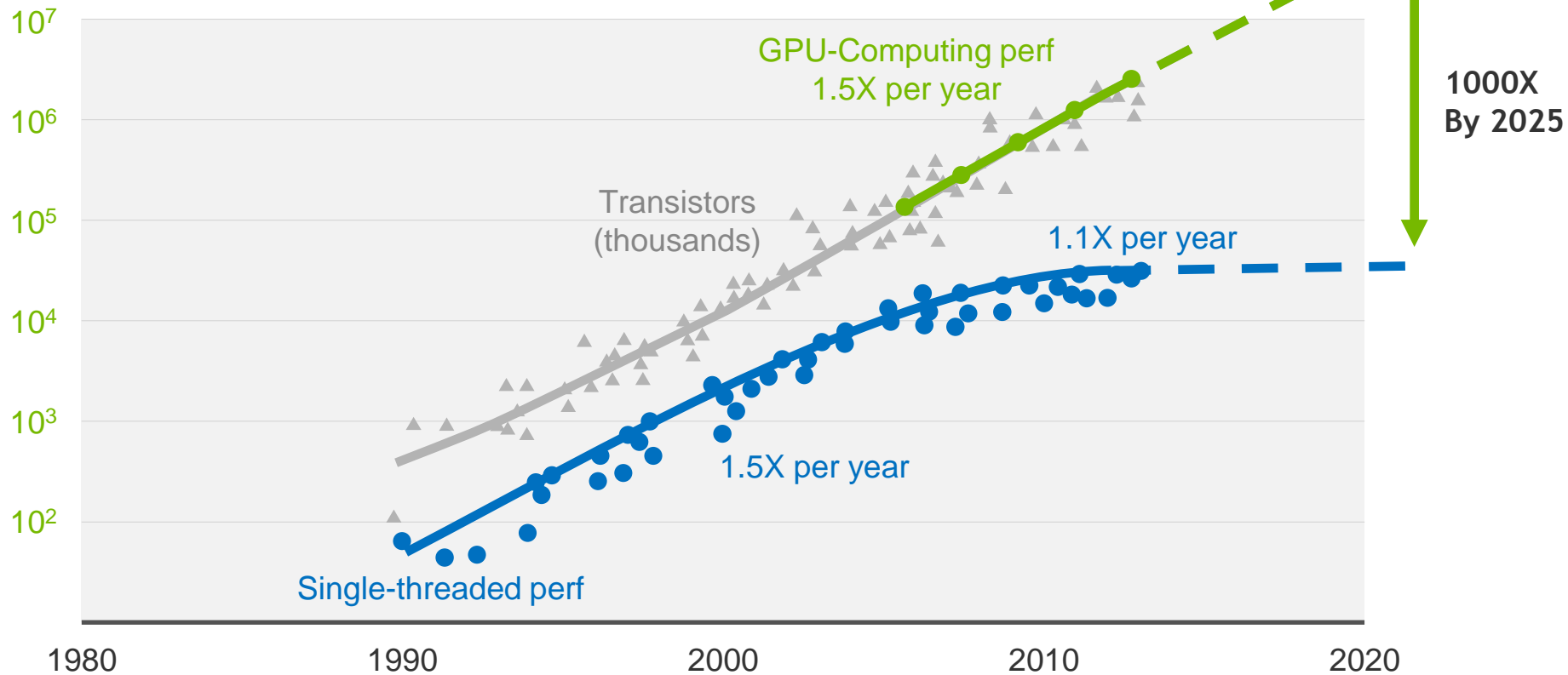
- Part 1: Gradient Descent
- Part 2: Stochastic Gradient Descent
- Part 3: Optimizing training with batch size

# CONTEXT: WHY USE MULTIPLE GPUS?

The background of the slide is a solid green color. Overlaid on this is a white network graph pattern. The graph consists of numerous small circular nodes connected by thin white lines, forming a complex, interconnected web that is denser on the right side and fades towards the left.

# TRENDS IN COMPUTATIONAL POWER

Historically we never had large datasets or compute



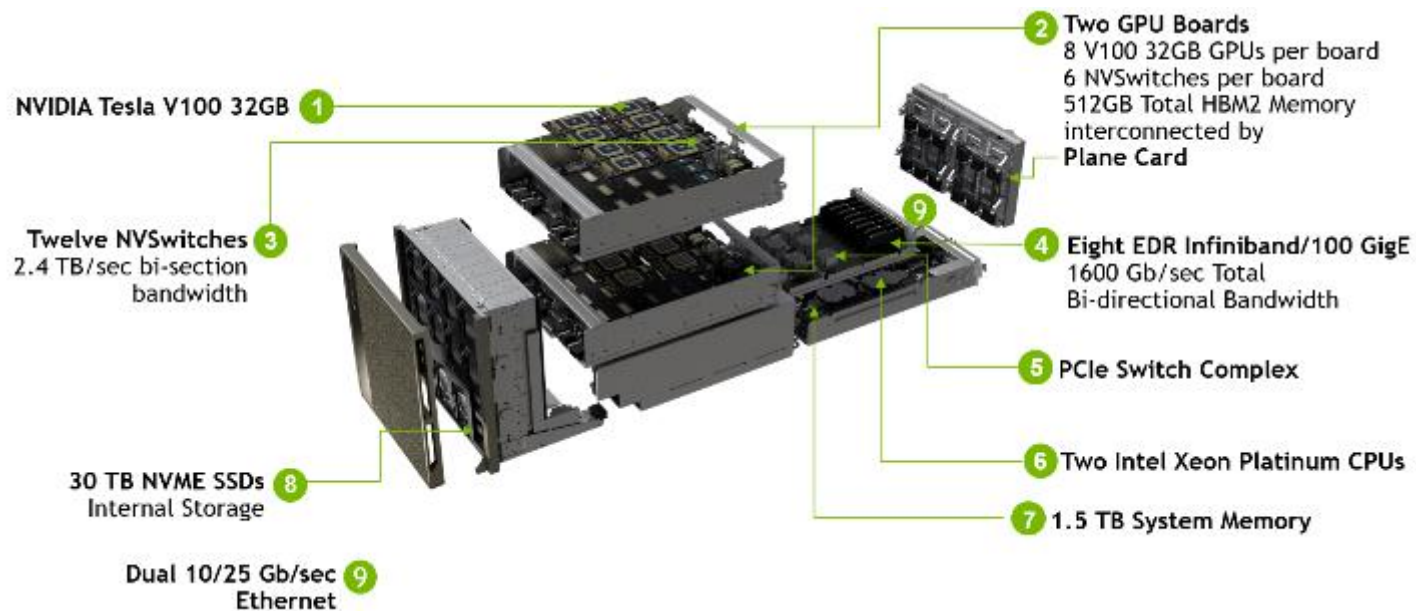
# TRENDS IN COMPUTATIONAL POWER

2 PF/s in November 2009



# TRENDS IN COMPUTATIONAL POWER

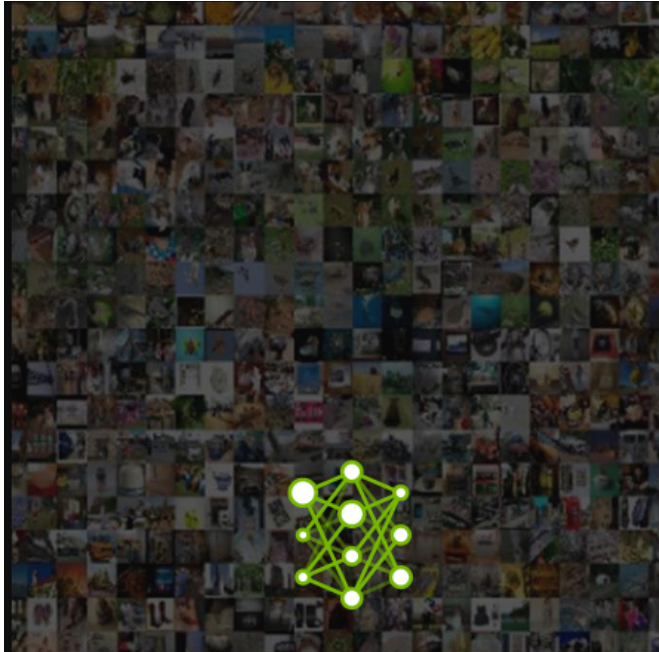
2 PF/s today



# NEURAL NETWORK COMPLEXITY IS EXPLODING

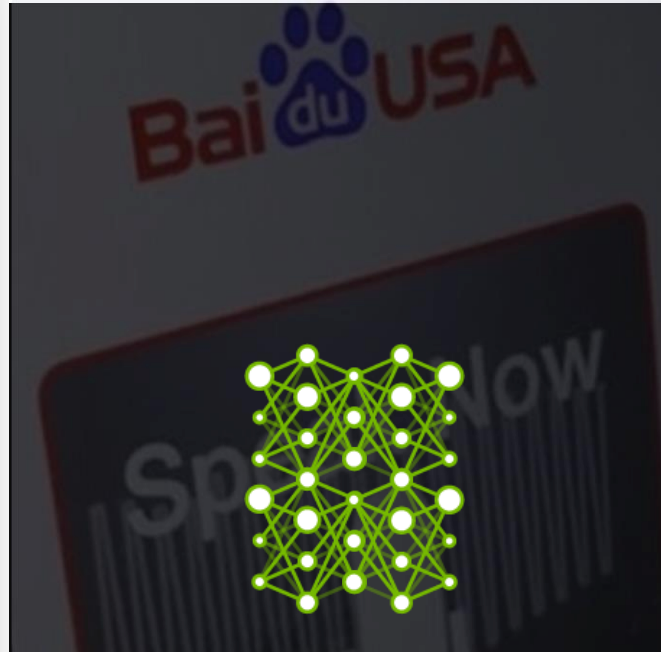
To Tackle Increasingly Complex Challenges

7 Exaflops  
60 Million Parameters



2015 - Microsoft ResNet  
Superhuman Image Recognition

20 Exaflops  
300 Million Parameters



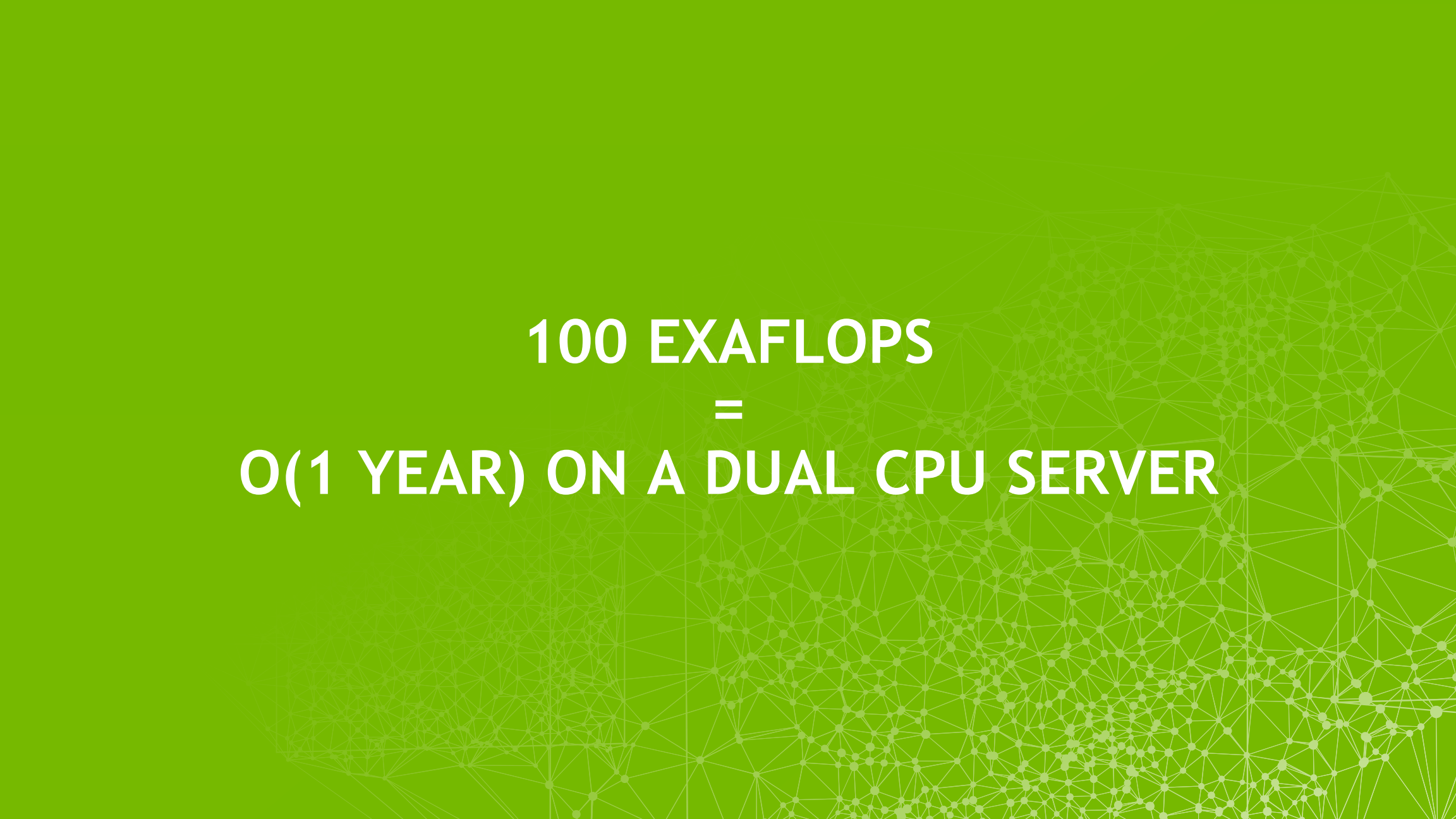
2016 - Baidu Deep Speech 2  
Superhuman Voice Recognition

100 Exaflops  
8700 Million Parameters



2017 - Google Neural Machine Translation  
Near Human Language Translation

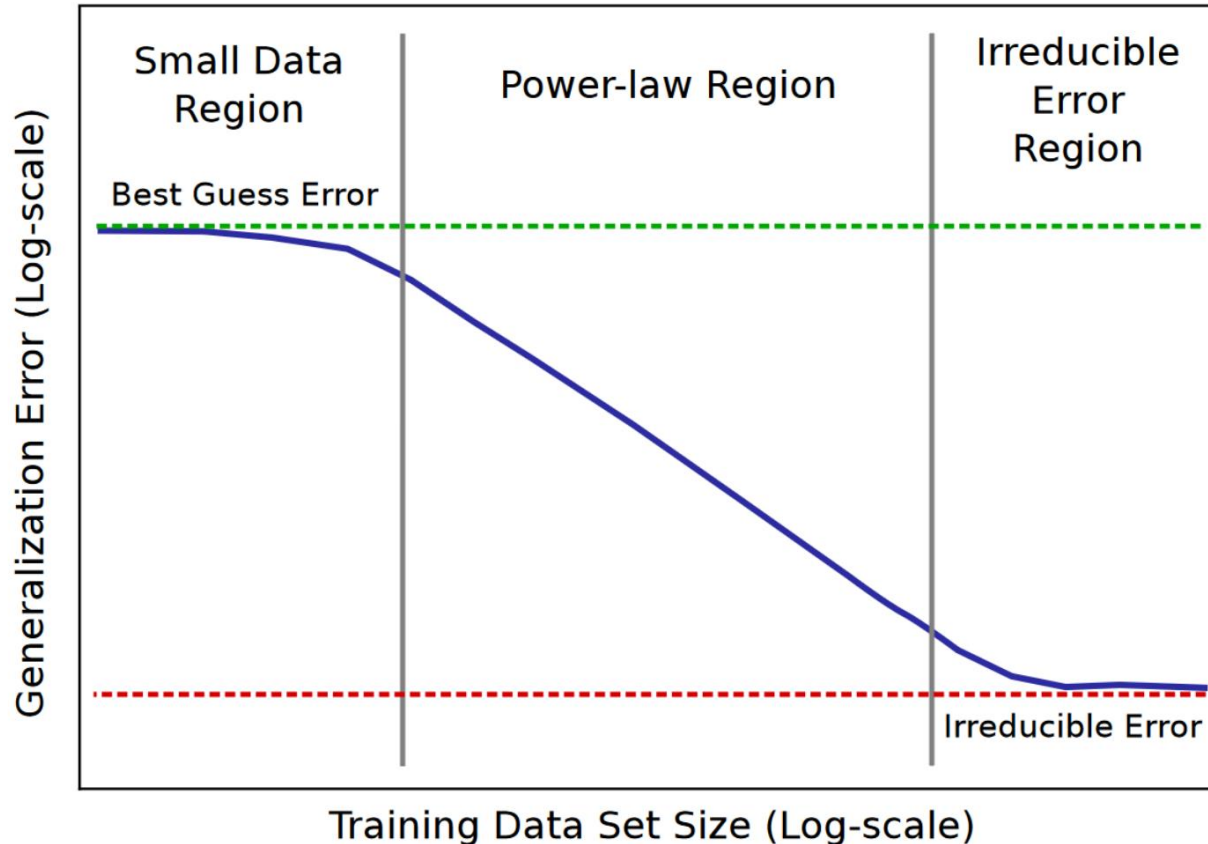




**100 EXAFLOPS  
=  
O(1 YEAR) ON A DUAL CPU SERVER**

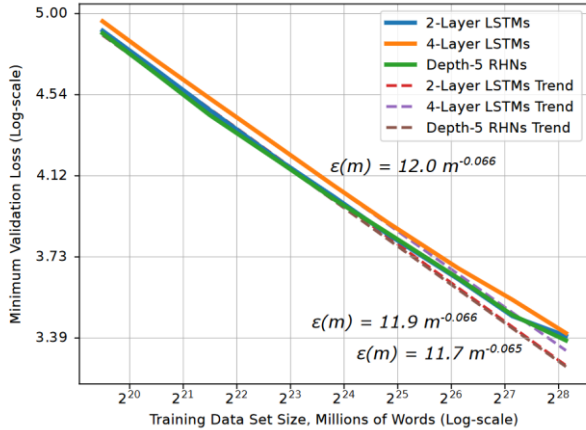
# EXPLODING DATASETS

Power-law relationship between dataset size and accuracy

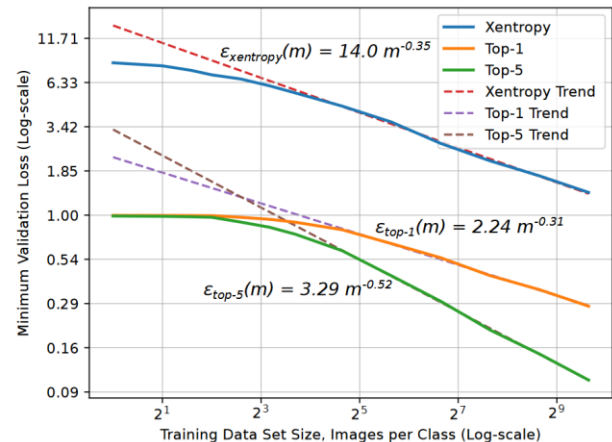
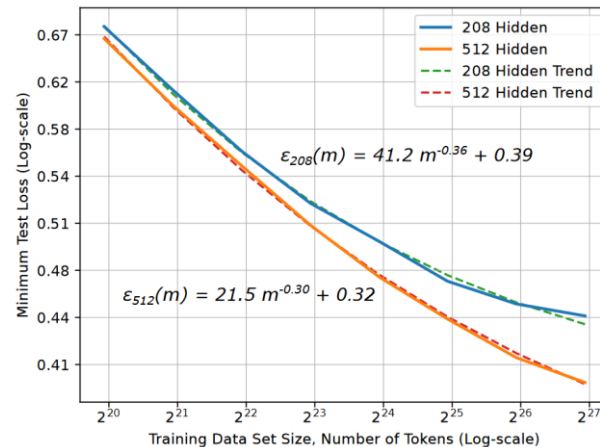
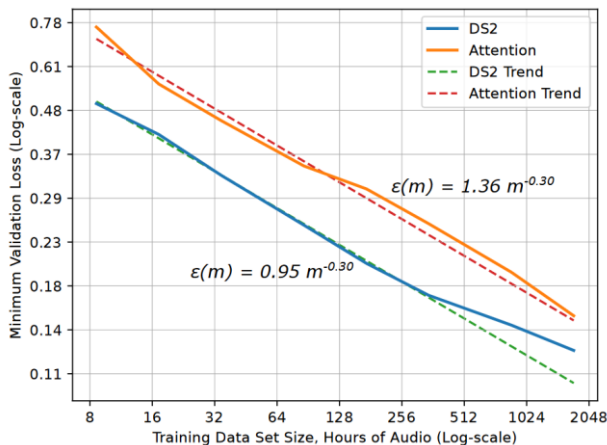
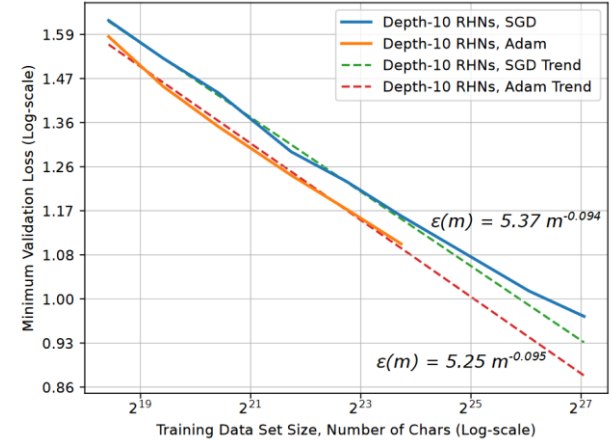


# EXPLODING DATASETS

Power-law relationship between dataset size and accuracy

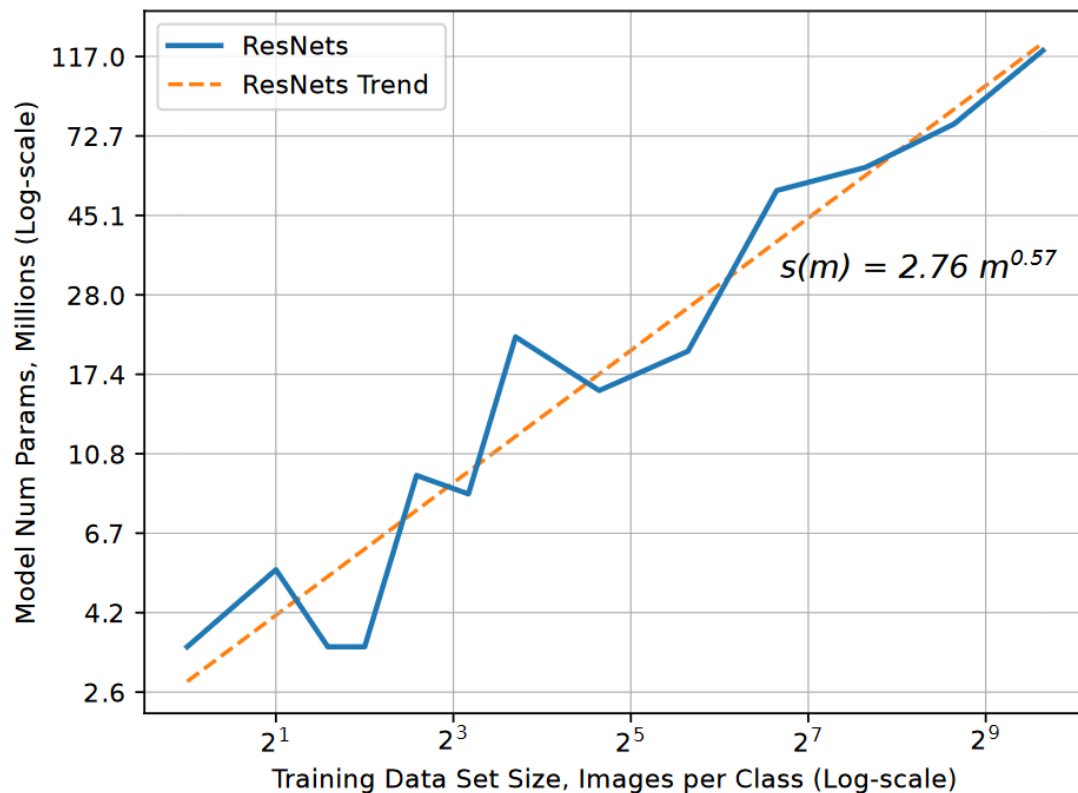
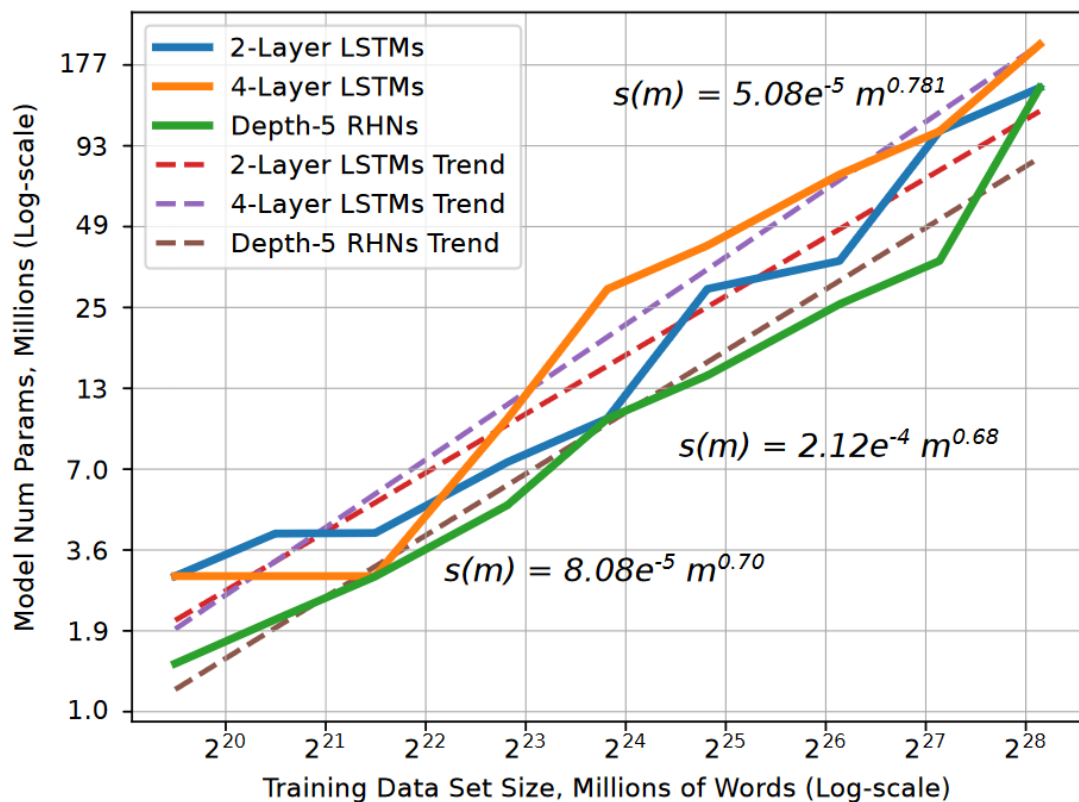


- Translation
- Language Models
- Character Language Models
- Image Classification
- Attention Speech Models



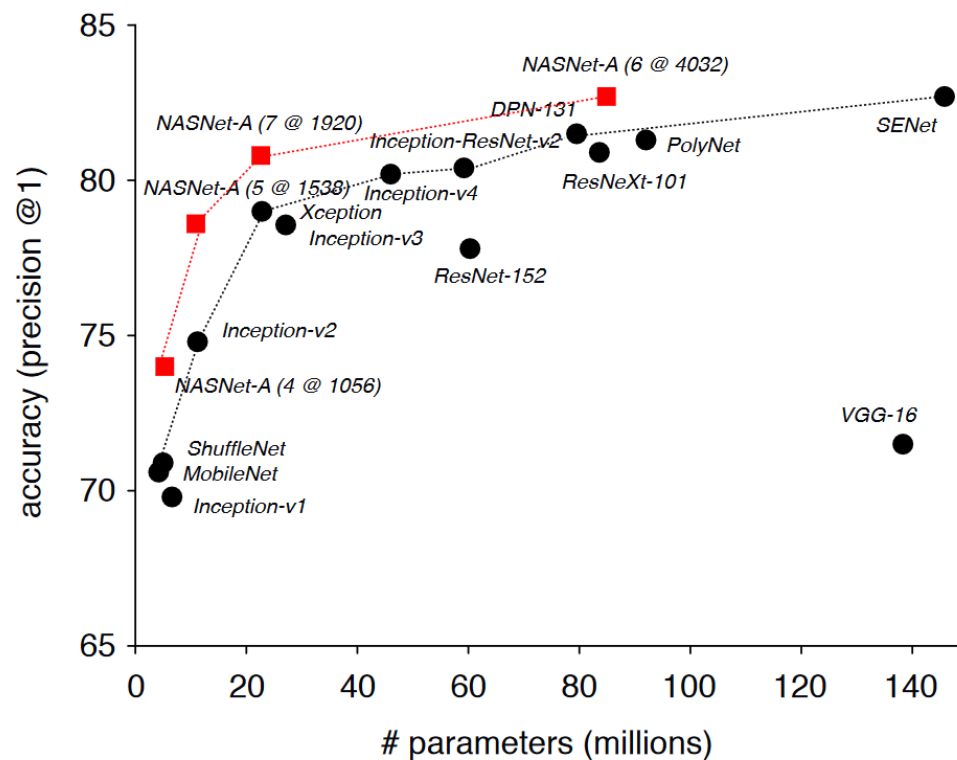
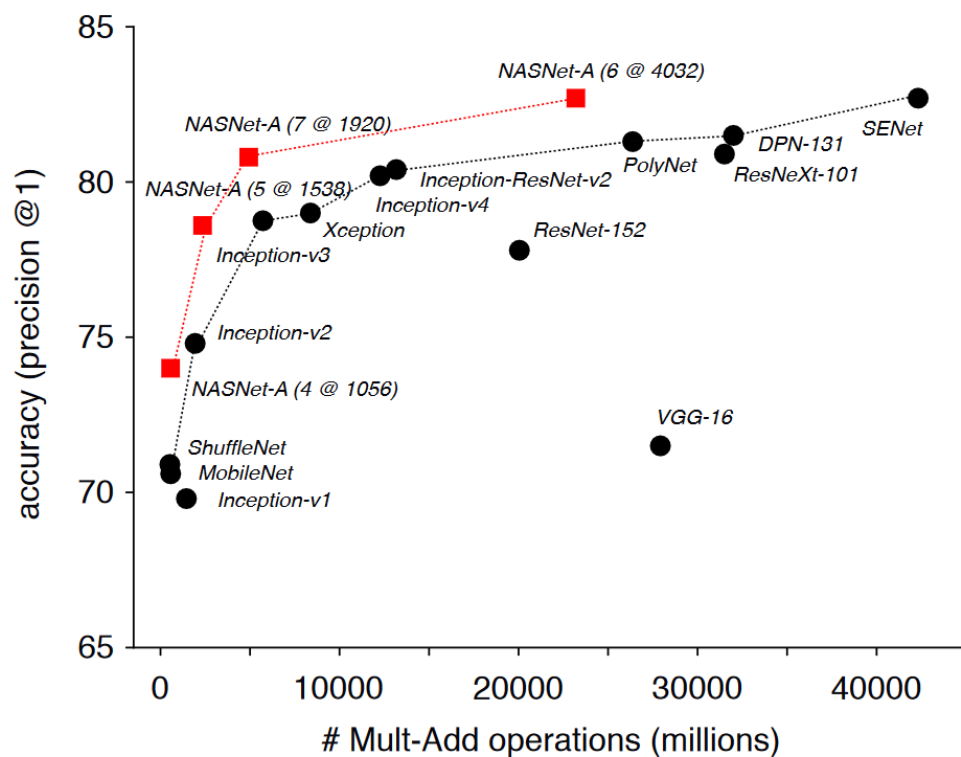
# EXPLODING MODEL COMPLEXITY

Though model size scales sublinearly



# EXPLODING MODEL COMPLEXITY

Though model size scales sublinearly



# IMPLICATIONS



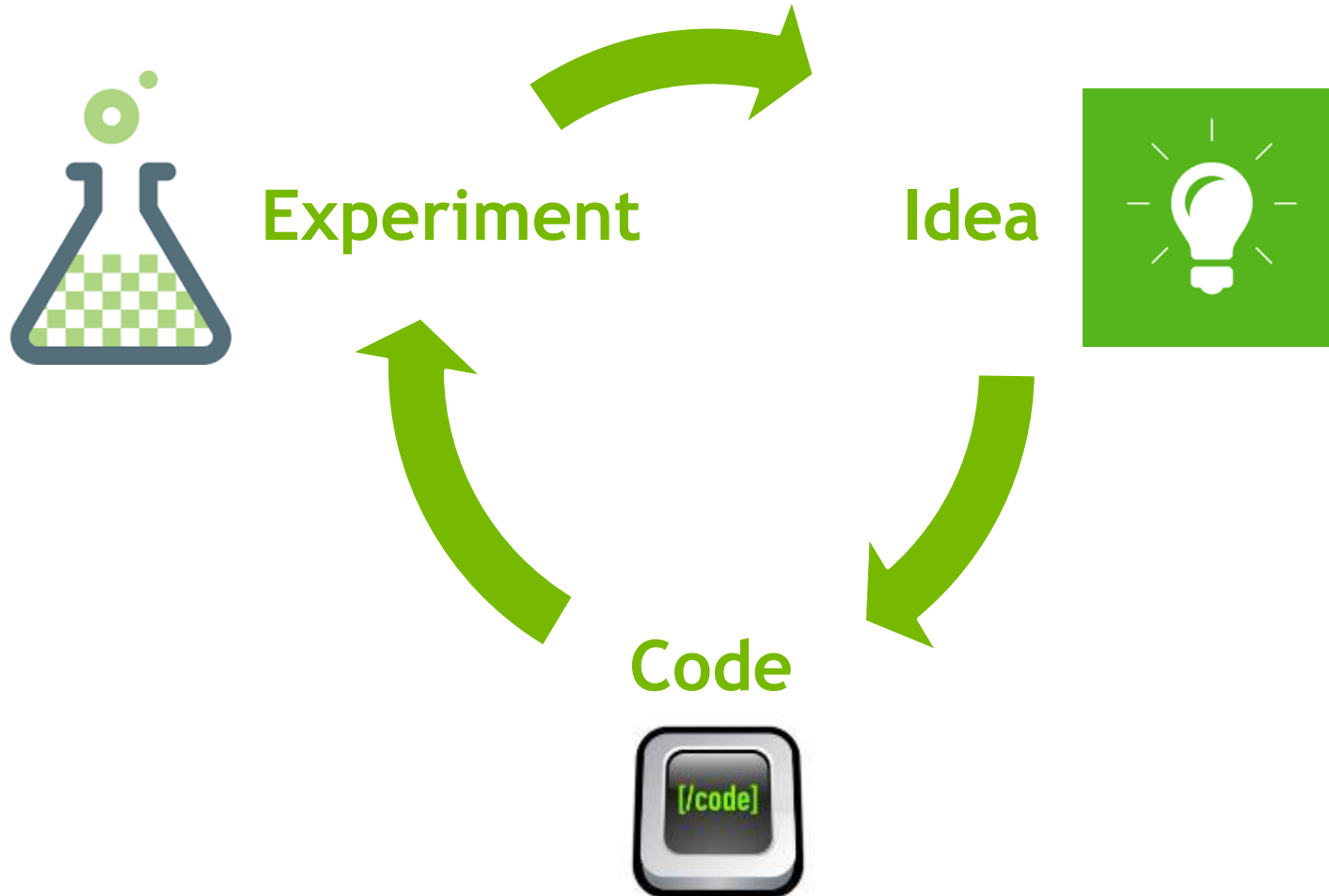
# IMPLICATIONS

## Good and bad news

- ▶ The good news: Requirements are predictable.
  - ▶ We can predict how much data we will need.
  - ▶ We can predict how much computing power we will need.
- ▶ The bad news: The values can be significant.
  - ▶ The silver lining is that deep learning has taken impossible problems and made them merely expensive.

# IMPLICATIONS

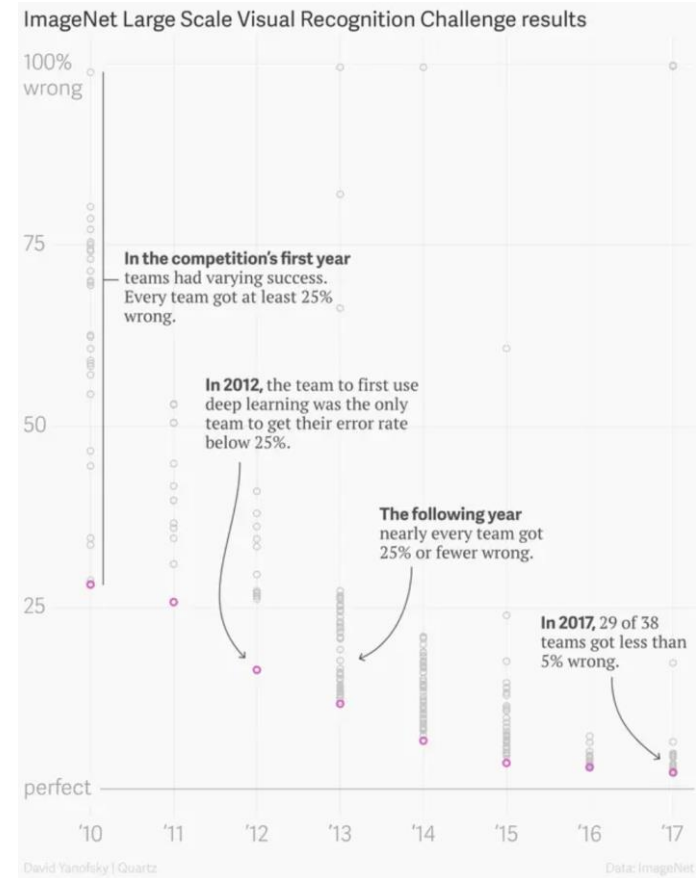
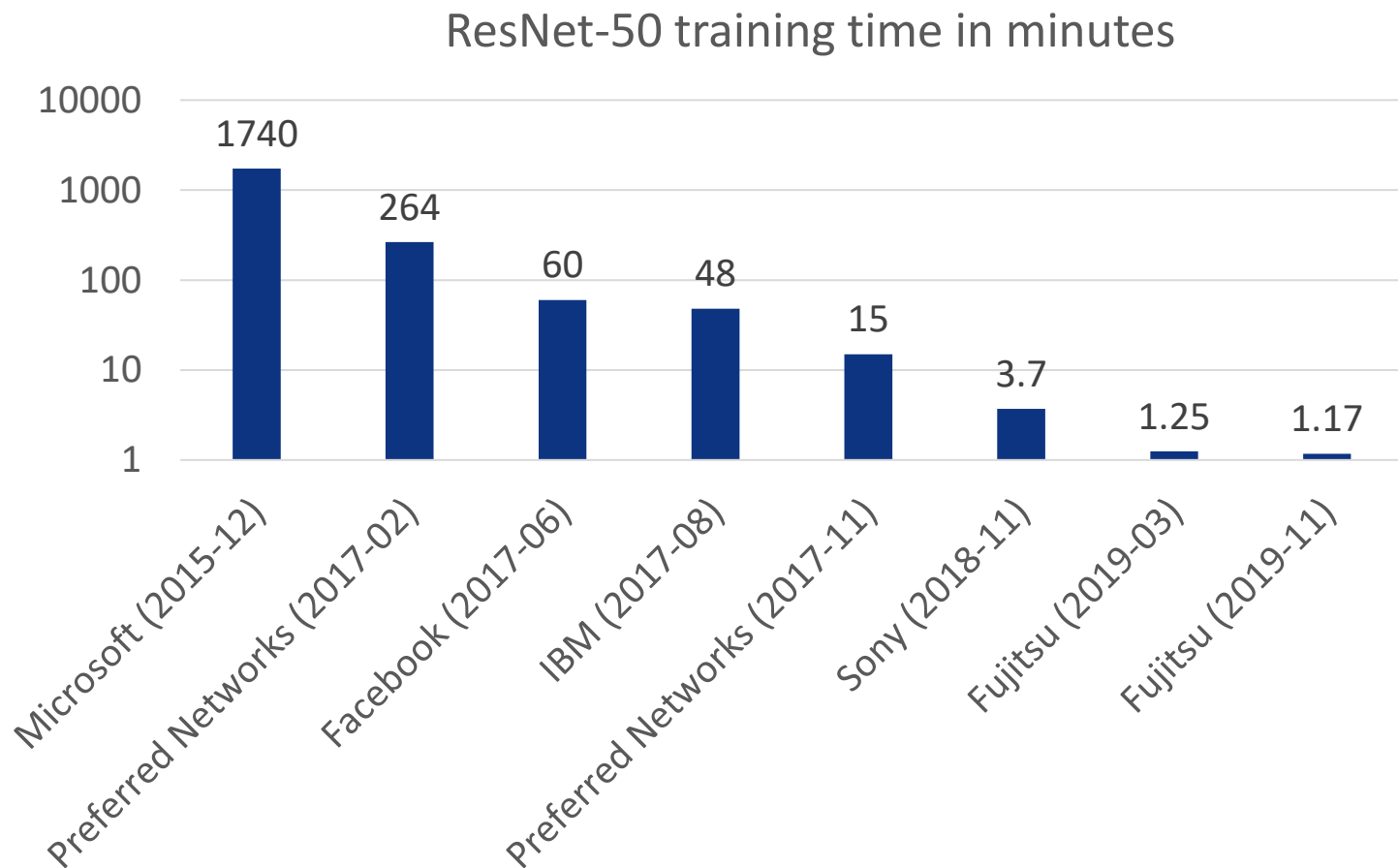
Deep learning is experimental; we need to train quickly to iterate





# ITERATION TIME

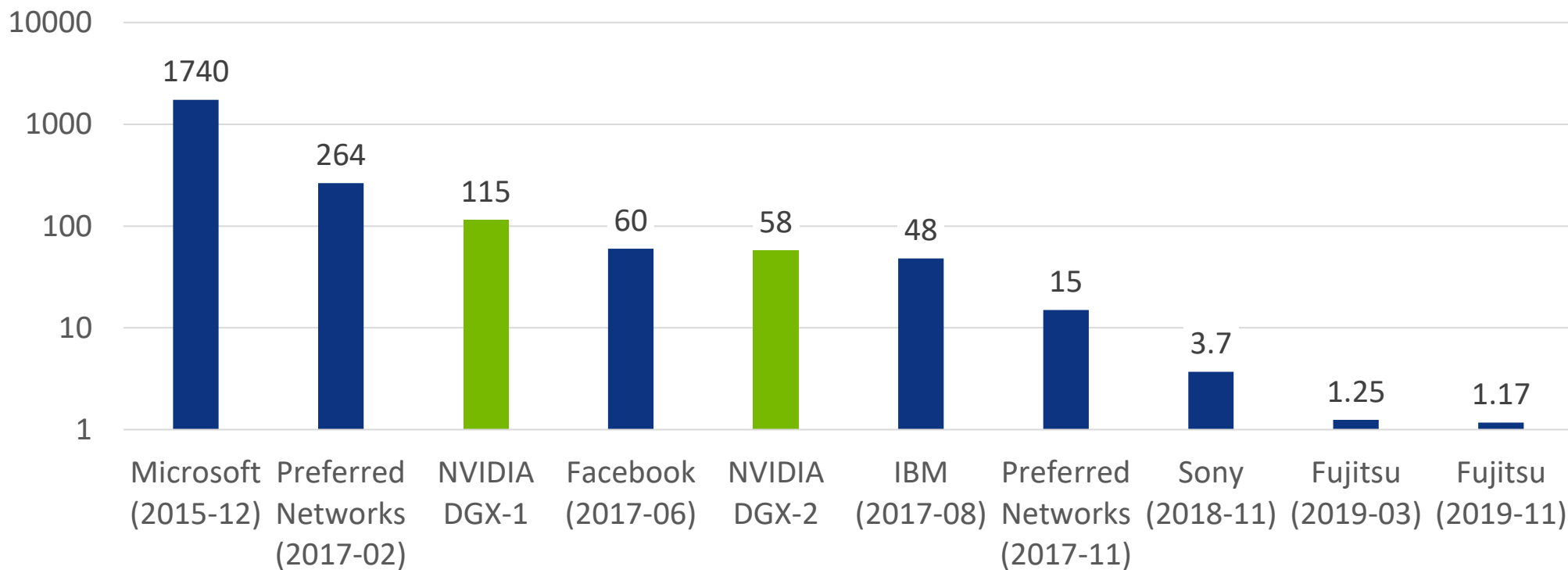
Short iteration time is fundamental for success



# ITERATION TIME

Short iteration time is fundamental for success

ResNet-50 training time in minutes

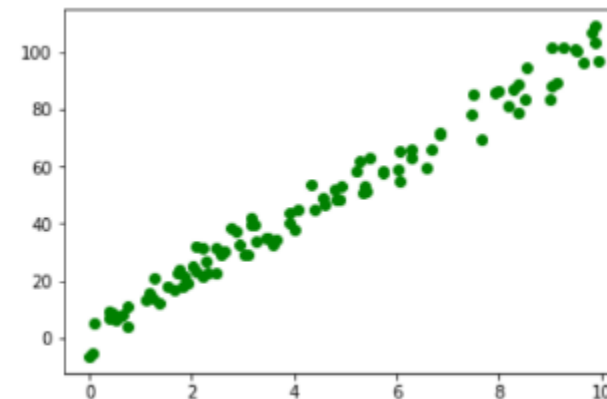
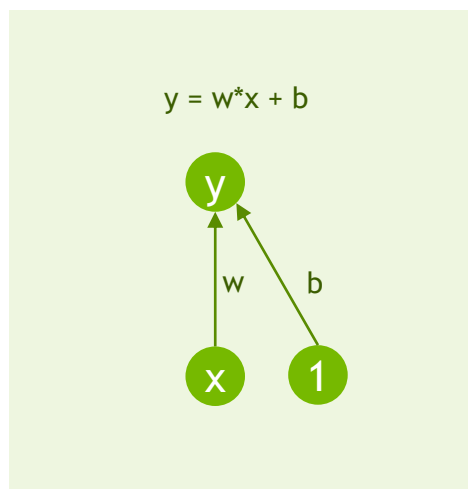


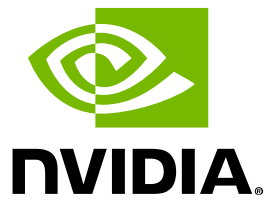


# INTRO TO THE LAB

# STARTING WITH A LINEAR MODEL

Our goal is to find best model parameters (combination of  $w$  and  $b$ ) to fit the data





DEEP  
LEARNING  
INSTITUTE

[www.nvidia.com/dli](http://www.nvidia.com/dli)