# AGENDA

**Deep Learning at Scale**
- Why scaling
- Tips and Tricks

**Desired outcomes**
- Understand when to and how to scale
- Know the typical techniques to apply
- Understand the theory/concepts of typical scaling techniques

# Deep Learning At Scale

**Proof-of-concept**

- Model architecture search
- Hyperparameter Optimization
- Toy dataset/Toy models
- Low frequency of retraining
- Non-optimized resource utilization

**Production**

- Models need to be trained and retrained with shorter times and higher frequency.
- Massive datasets
- Big & complex models

# Pain Points



DL models training:

- o **Time consuming**: can take days, weeks...
- o **Capability:** is limited by memory capacity on batch size and model size

Scaling is necessary but hard:

- o Convergence and Stability
- o Computation and Scaling efficiency
- o Hardware Limits

# Techniques for Scaling Deep Learning Training

- **Convergence and Stability**
  - Warmup
  - Linear Scaling Rule
  - LARS

- **Computation and Scaling efficiency**
  - Automatic Mixed Precision

- **Hardware limits on Dataset and Model Size**
  - User profiles/preference of millions of users
  - Data Parallelism
  - Model Parallelism

# Scaling Success

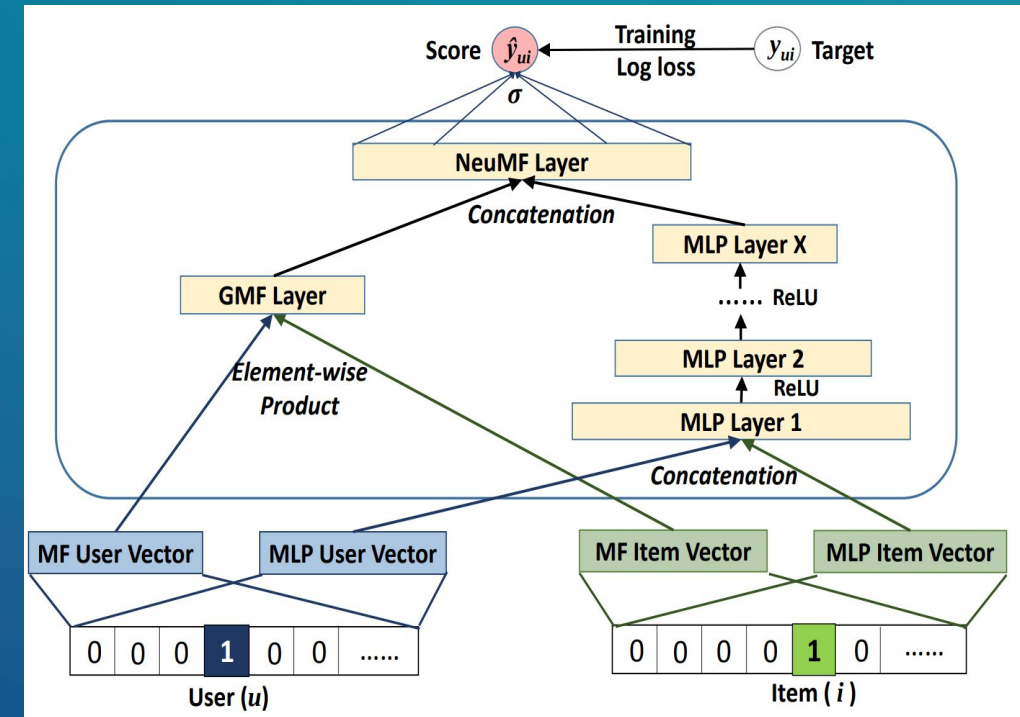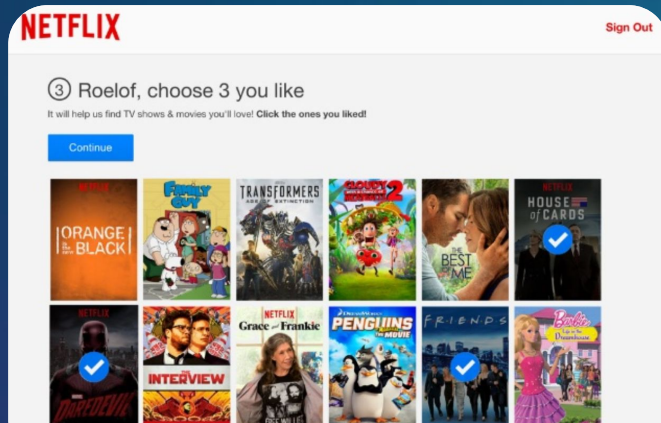| BERT PreTraining on DGX SuperPOD | |
| --- | --- |
| # V100 GPUs | Time to train (hours) |
| 16 | 58.4 (2.4 days) |
| 256 | 3.9 |
| 1024 | 1.2 |
| 1472 | 0.9 (53 min) |

**GOAL :**
**1. Maintain Accuracy**
**2. Decrease Time to Train**

**Let's do it on a smaller model: NCF BERT has 110 million parameters!**

# Recommender System: Neural Collaborative Filtering

**Recommendation engines are everywhere...**
- **Personalized**
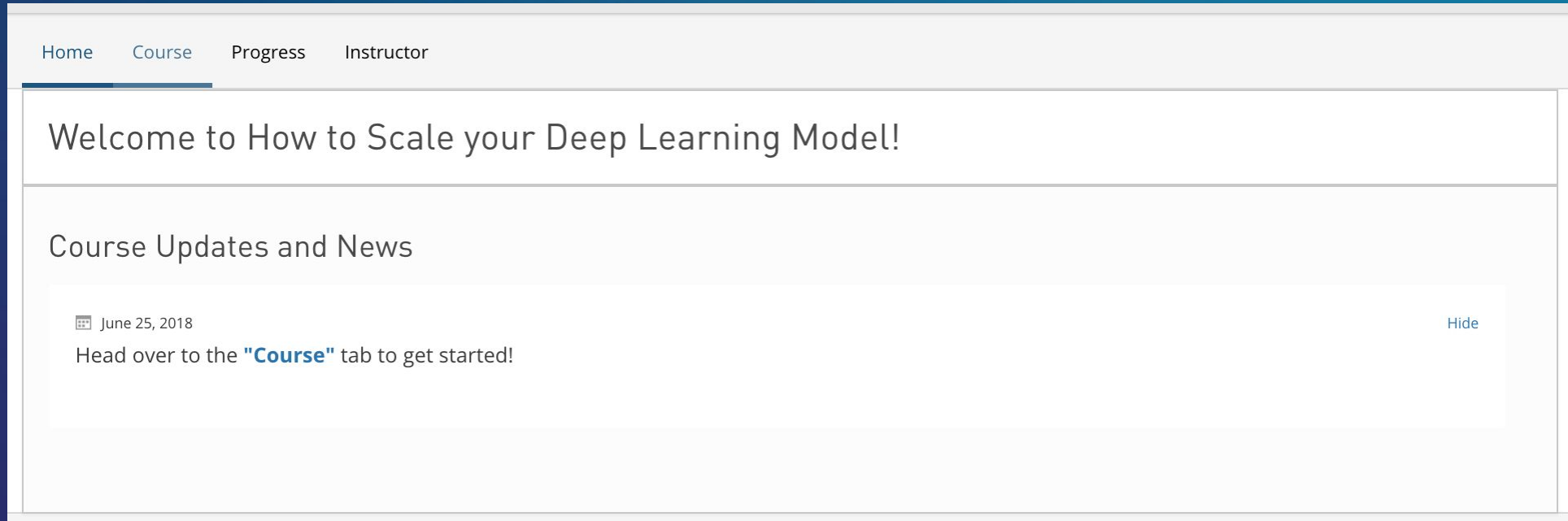- **Convenient**
- **More efficient**





**Deep Learning based Recommendation system model architecture**

# Launch Hands-on Task

- Create Account at https://courses.nvidia.com/join
- Go to courses.nvidia.com/dli-event
- Browser Recommendation: Chrome
- Use event code and create an account

Select the "Course" Tab

Home    Course    Progress    Instructor

## Welcome to How to Scale your Deep Learning Model!

### Course Updates and News

📅 June 25, 2018                                                                          Hide

Head over to the **"Course"** tab to get started!

# Open the first hands-on section

How to Scale your Deep Learning Model

Search th

How to Scale your Deep Learning Model

Click here to get started                                    Resume Course ➔

Feedback

# Select the Start button and wait

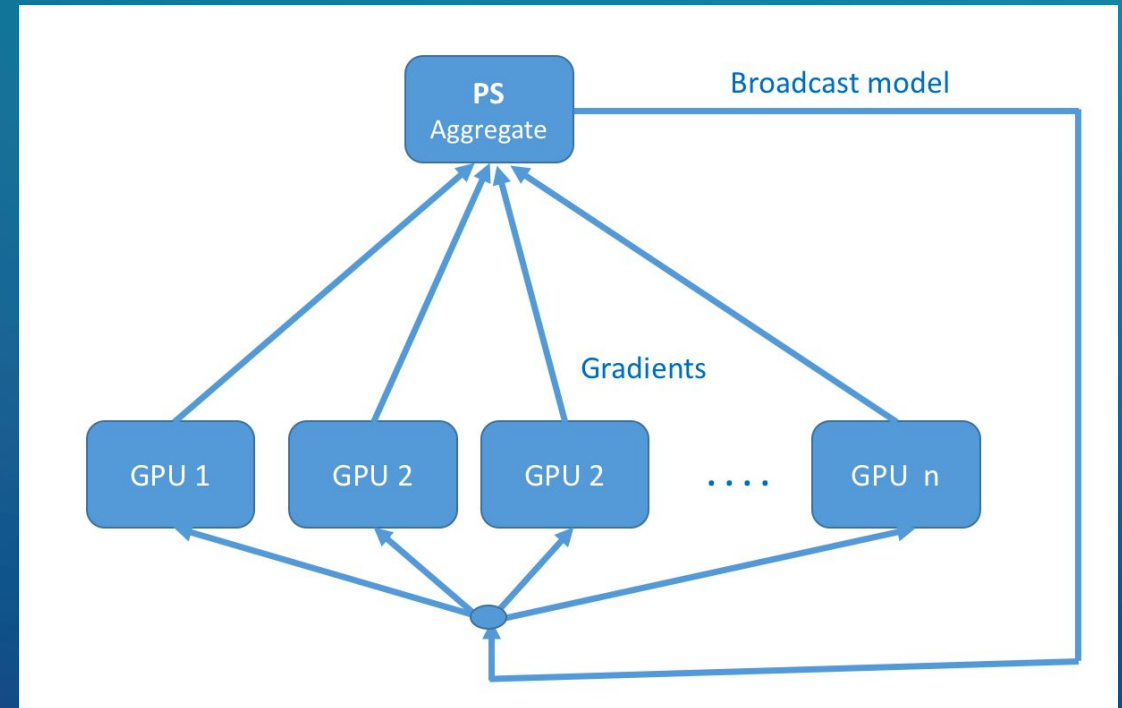Launch Lab                                                          VIEW UNIT IN STUDIO

🔖 Bookmark this page

NVIDIA. | DEEP LEARNING INSTITUTE                              ▶ START

Please click the "Start" button to get started.

#GHC19

# What else can we do?

- Is 192 still the maximum batch size we can use?
- Can we scale further than 1 GPU?
  - Data Parallelism

https://github.com/NVIDIA/DeepLearningExamples



*Image from Towards Data Science

# Summary

Productivity matters : teams with better tools/scaling can try out more ideas
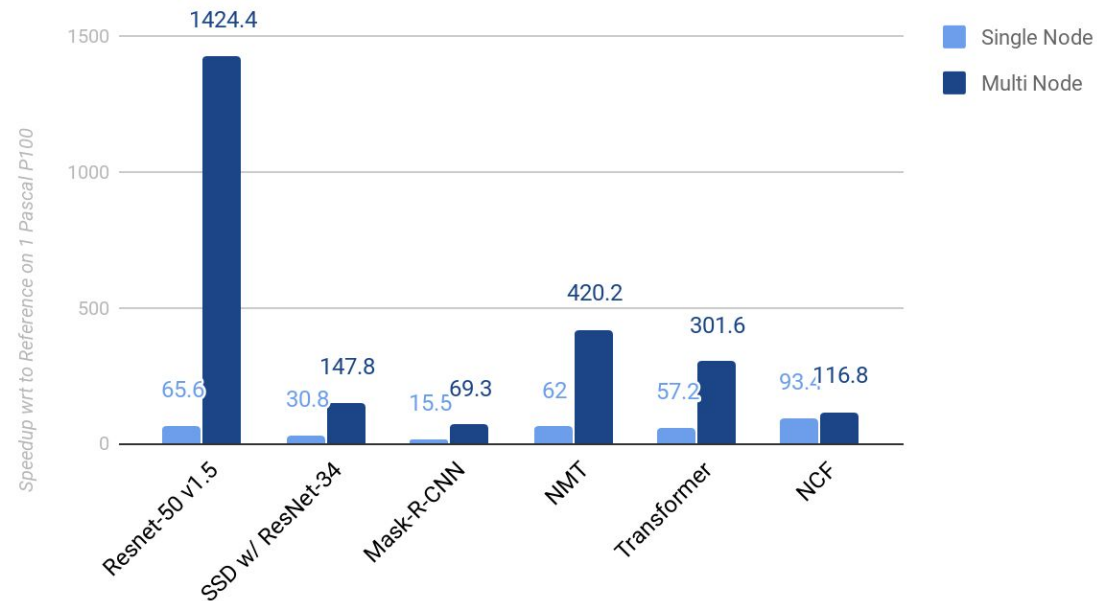
| SGD<br>BS = 4096<br>~1230 sec | → | + LR Scaling<br>BS *= 16<br>~140 sec | → | + WARMUP<br>BS *= 192<br>~100 sec | → | +LARS<br>BS *= 192<br>~110 sec | → | +AMP<br>BS *= 192<br>~ 55 sec |



MLPerf v0.5 Results