



Deep Learning in Parallel

Alan Sussman, Department of Computer Science



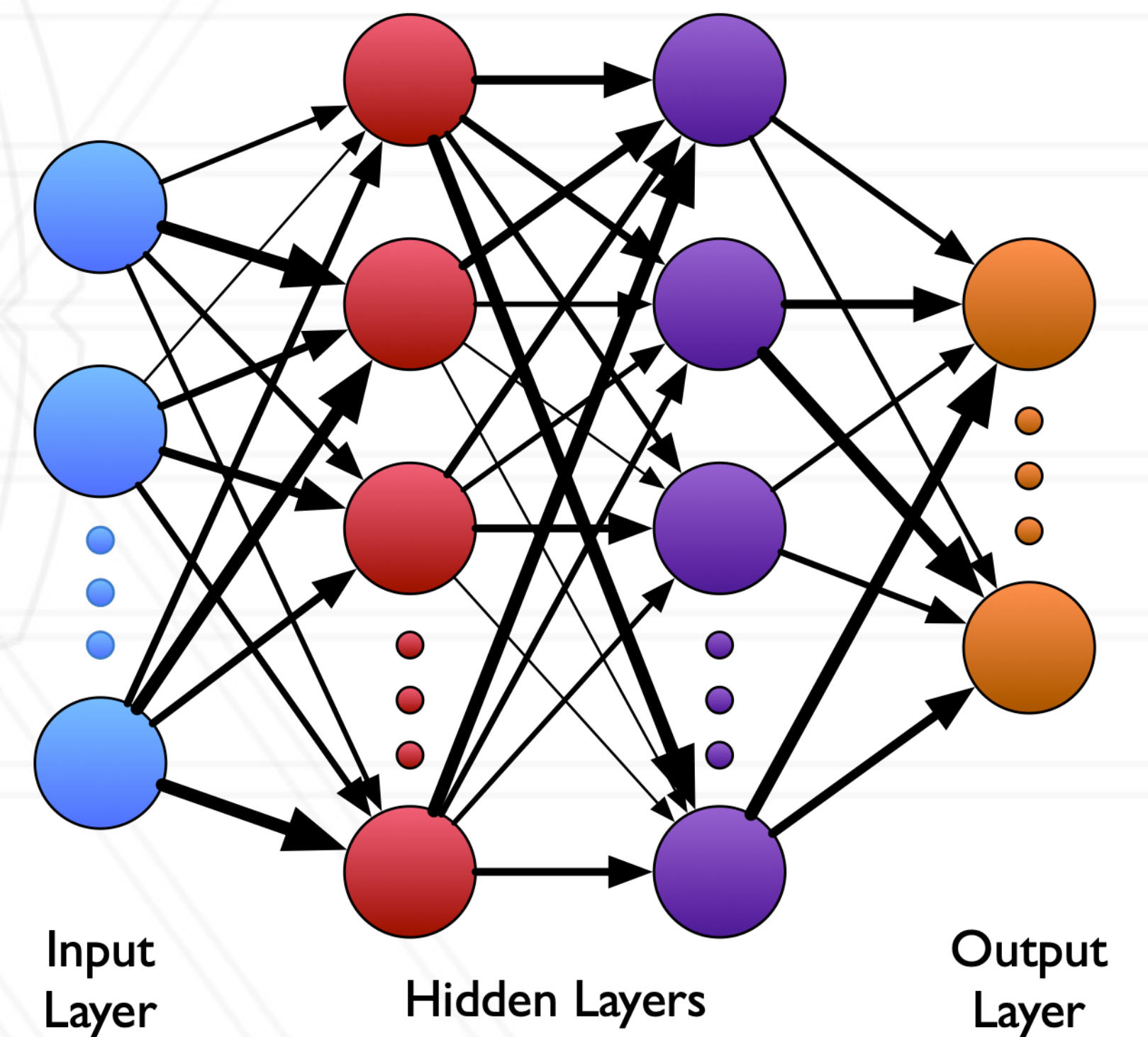
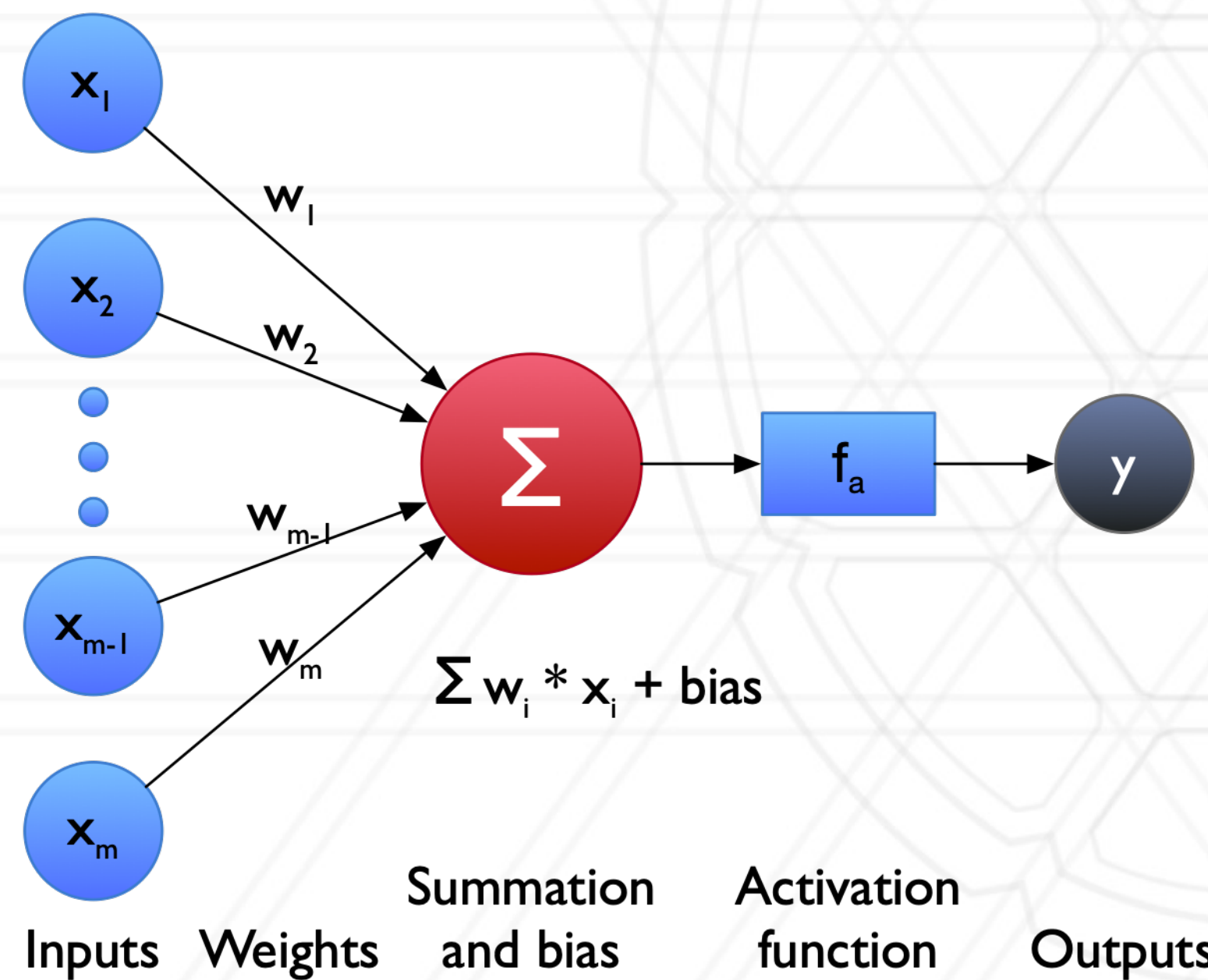
UNIVERSITY OF
MARYLAND

Announcements

- Quiz 3 will be posted on Wednesday, May 10, 11AM
 - In ELMS, for 24 hours
 - Mainly on topics since last quiz
- Course evaluation: <https://www.courseevalum.umd.edu>

Deep neural networks

- Neural networks can be used to model complex functions
- Several layers that process “batches” of the input data

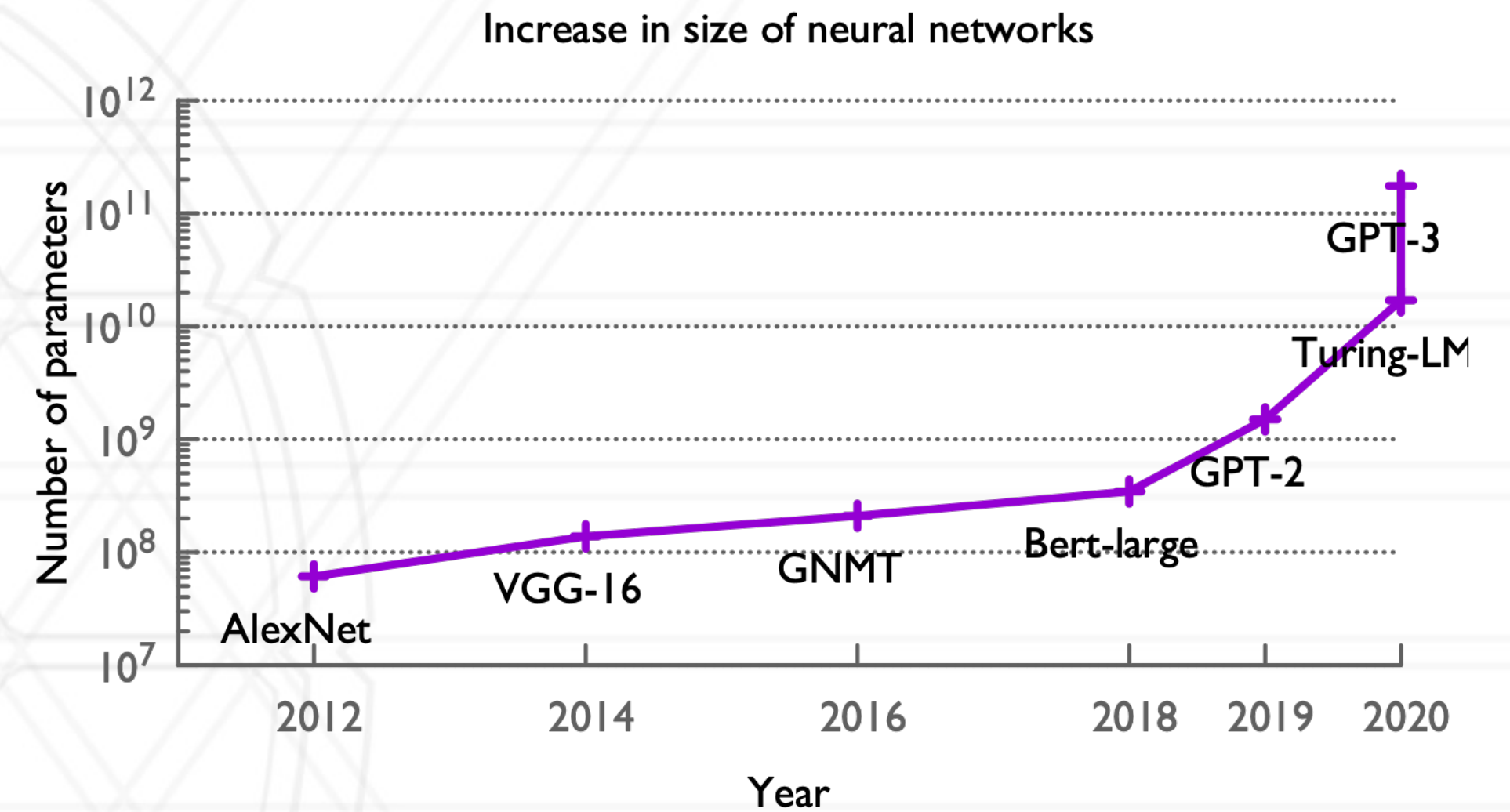


Other definitions

- Learning/training: task of selecting weights that lead to an accurate function
- Loss: a scalar proxy that when minimized leads to higher accuracy
- Gradient descent: process of updating the weights using gradients (derivatives) of the loss weighted by a learning rate
- Mini-batch: Small subsets of the dataset processed iteratively
- Epoch: One pass over all the mini-batches

Parallel/distributed training

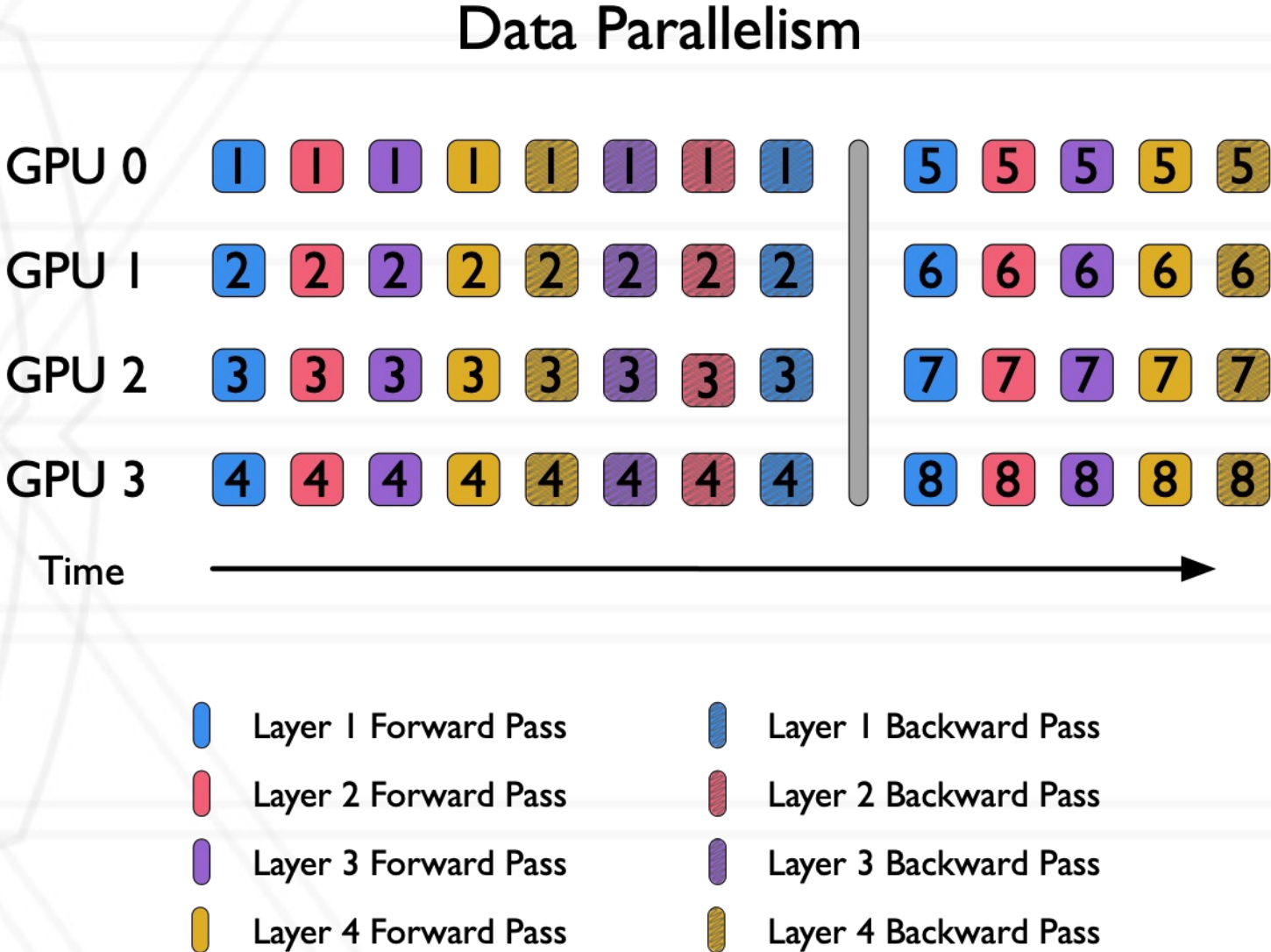
- Many opportunities for exploiting parallelism
- Iterative process of training (epochs)
- Many iterations per epoch (mini-batches)
- Many layers in DNNs



Framework	Type of Parallelism	Largest Accelerator Count	Largest Trained Network (No. of Parameters)
FlexFlow	Hybrid	64 GPUs	24M*
PipeDream	Inter-Layer	16 GPUs	138M
DDP	Data	256 GPUs	345M
GPipe	Inter-Layer	8 GPUs	557M
MeshTensorFlow	Intra-Layer	512-core TPUv2	4.9B
Megatron	Intra-Layer	512 GPUs	8.3B
TorchGPipe	Inter-Layer	8 GPUs	15.8B
KARMA	Data	2048 GPUs	17B
LBANN	Data	3072 CPUs	78.6B
ZeRO	Data	400 GPUs	100B

Data parallelism

- Divide training data among workers (GPUs)
- Each worker has a full copy of the entire NN and processes different mini-batches
- All-reduce operation to synchronize gradients

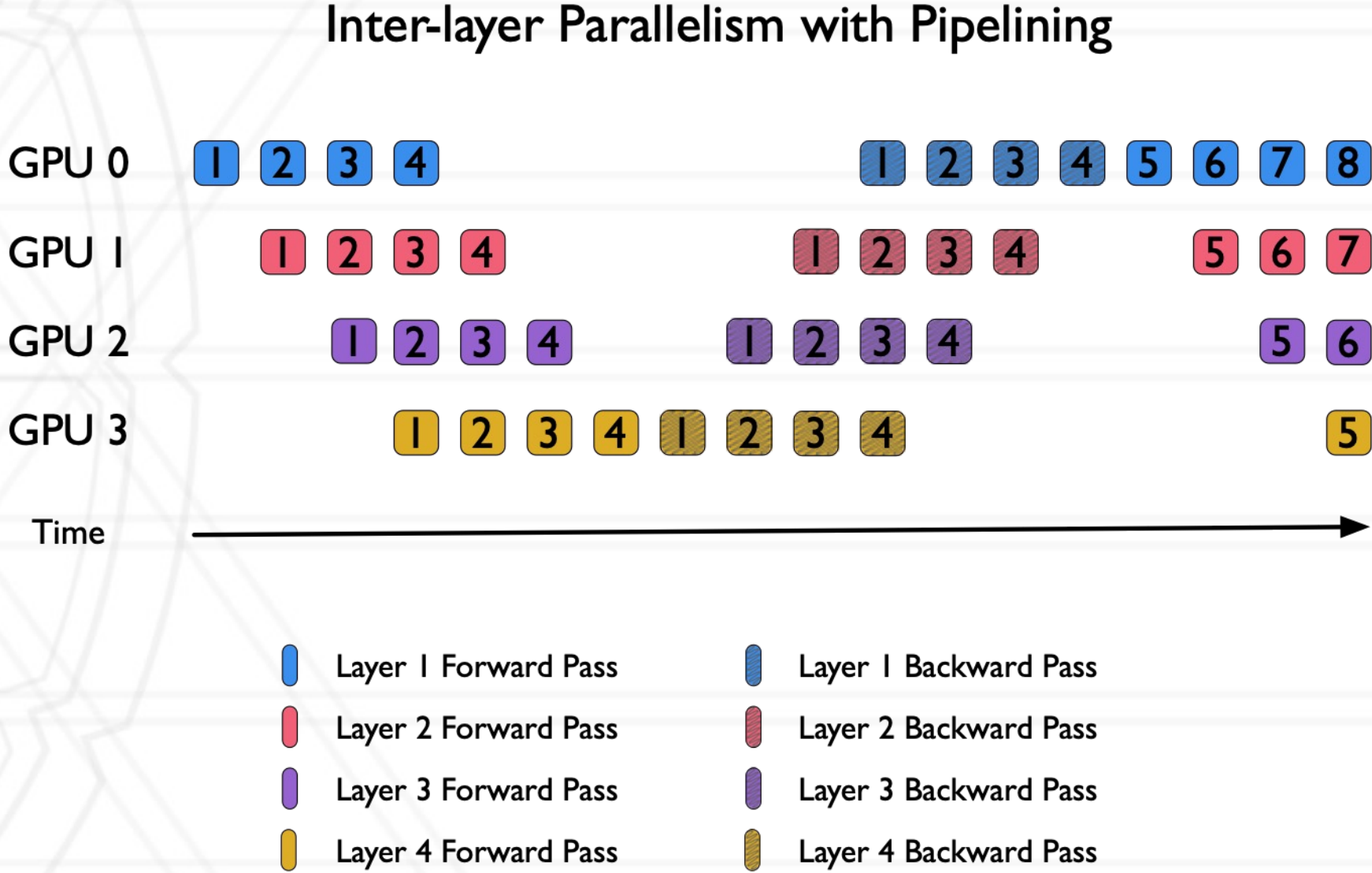


Intra-layer parallelism

- Enables training neural networks that would not fit in memory of a single GPU
- Distribute the work within a layer between multiple processes/GPUs

Inter-layer parallelism

- Distribute entire layers to different processes/GPUs
- Map contiguous subsets of layers
- Point-to-point communication (activations and gradients) between processes/GPUs managing different layers
- Use a pipeline of mini-batches to enable concurrent execution



Hybrid parallelism

- Using two or more approaches together in the same parallel framework
- 3D parallelism: use all three
- Popular serial frameworks: pytorch, tensorflow
- Popular parallel frameworks: DDP, MeshTensorFlow, Megatron-LM, ZeRO

Training vs. inference

- We talked about training, since that is very computationally intensive
- But once the DNN is trained, it is then used to do the ML task it has been designed to do (*inference*) – given an input (often not one that was in the input training set), produce the corresponding output
 - Classification
 - Pattern matching
 - ...
- Inference is much less computationally demanding than training, but will be done many times, potentially on edge devices (e.g., your smart phone)

Questions?



UNIVERSITY OF
MARYLAND