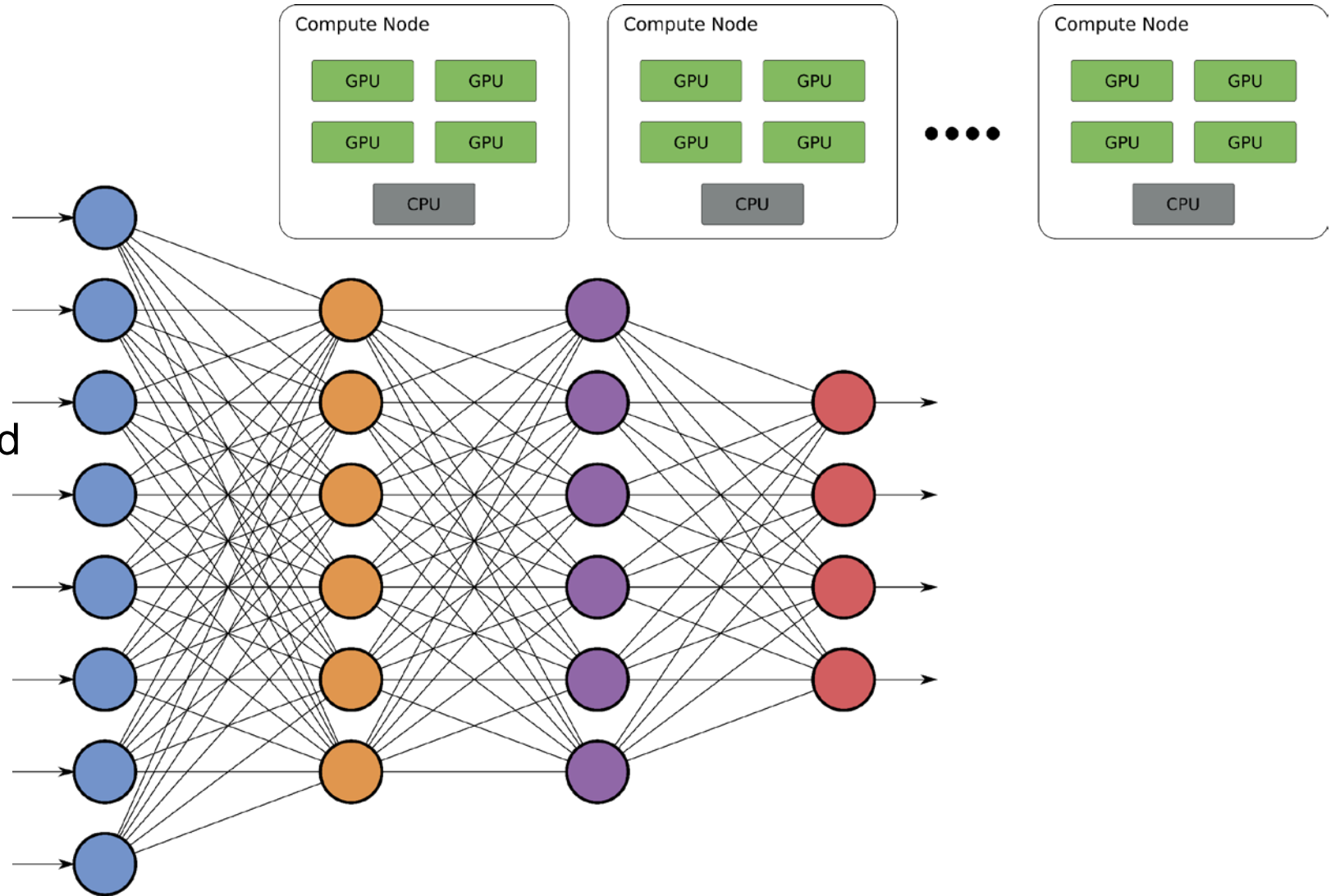


# Distributed Deep Learning

Feroz Zahid, Michael  
Riegler, and Tor Skeie  
Simula Research Laboratory and  
University of Oslo



- Full professor, Department of informatics, University of Oslo, Norway
- Adjunct research scientist, Simula Research Laboratory, Norway
- Fabriscale Technologies, CTO and co-founder
- Areas of expertise
  - HPC networking and management/middleware systems
  - Cloud computing and virtualization
  - Industrial Ethernet and wireless networking



[ **simula** . research laboratory ]



**This presentation will introduce distributed deep learning, walk through prominent techniques, and identify existing challenges and future directions**



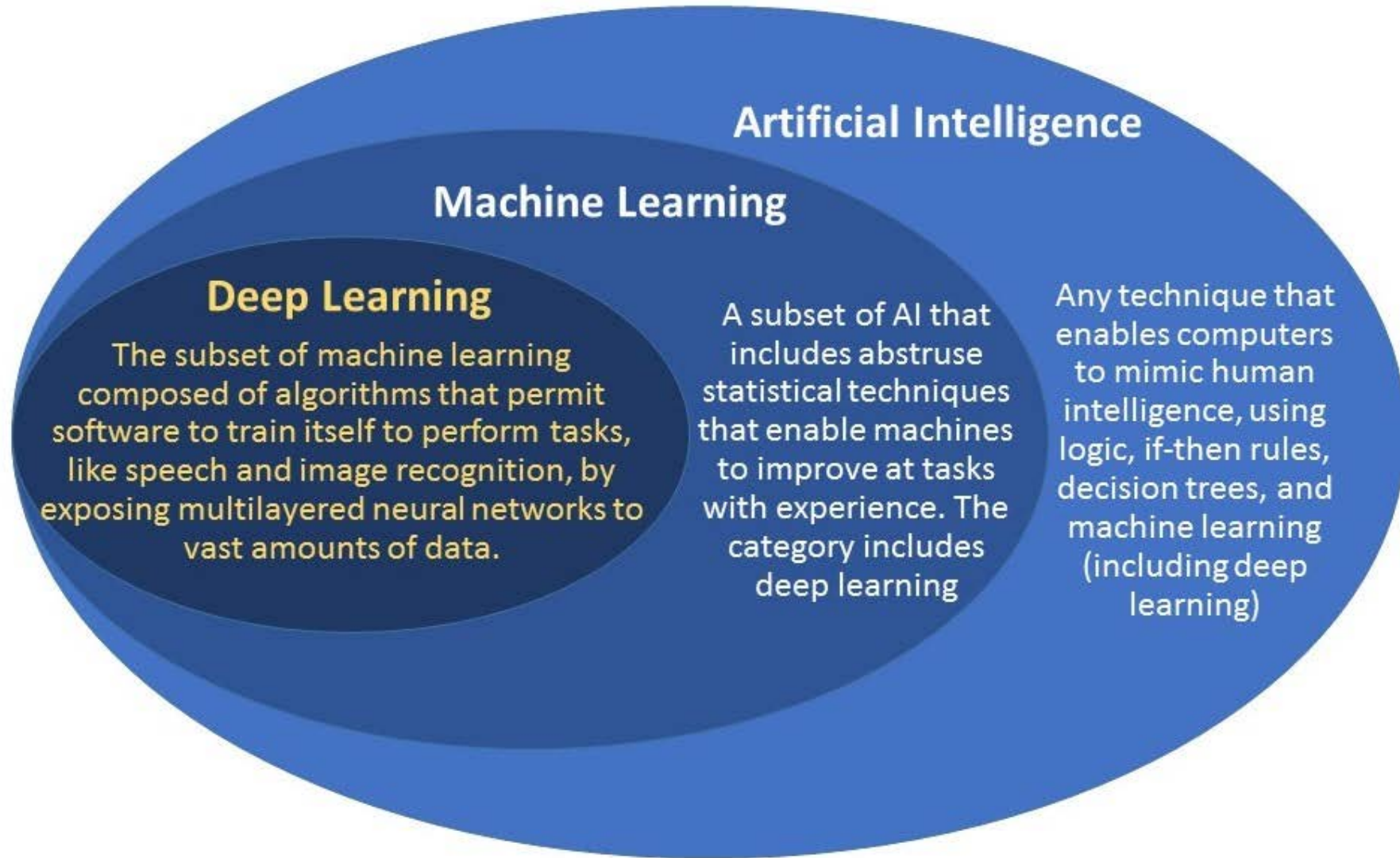
## **Introduction and Motivation**



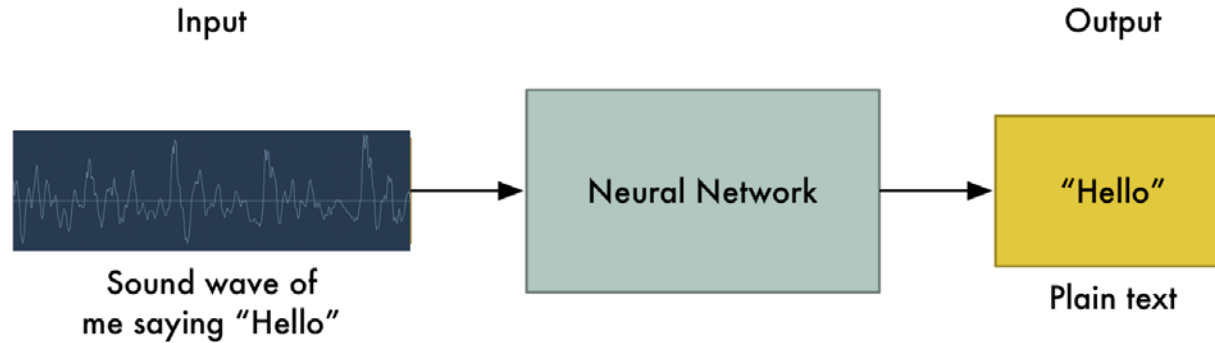
## **Existing Techniques and Toolsets**



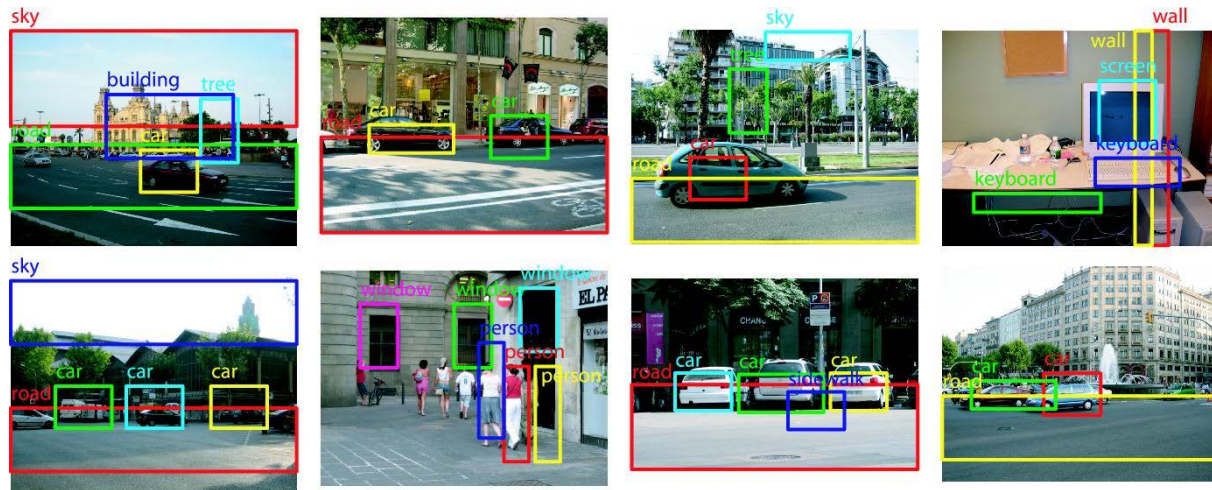
## **Future Directions**



# What you can do...



- \* Image classification
- \* Text processing
- \* Traffic classification
- \* Predictions of characteristics
- \* Climate prediction
- \* Health care
- \* many more....



Deep learning in Computer Vision, Fei-Fei Li et al., 2016

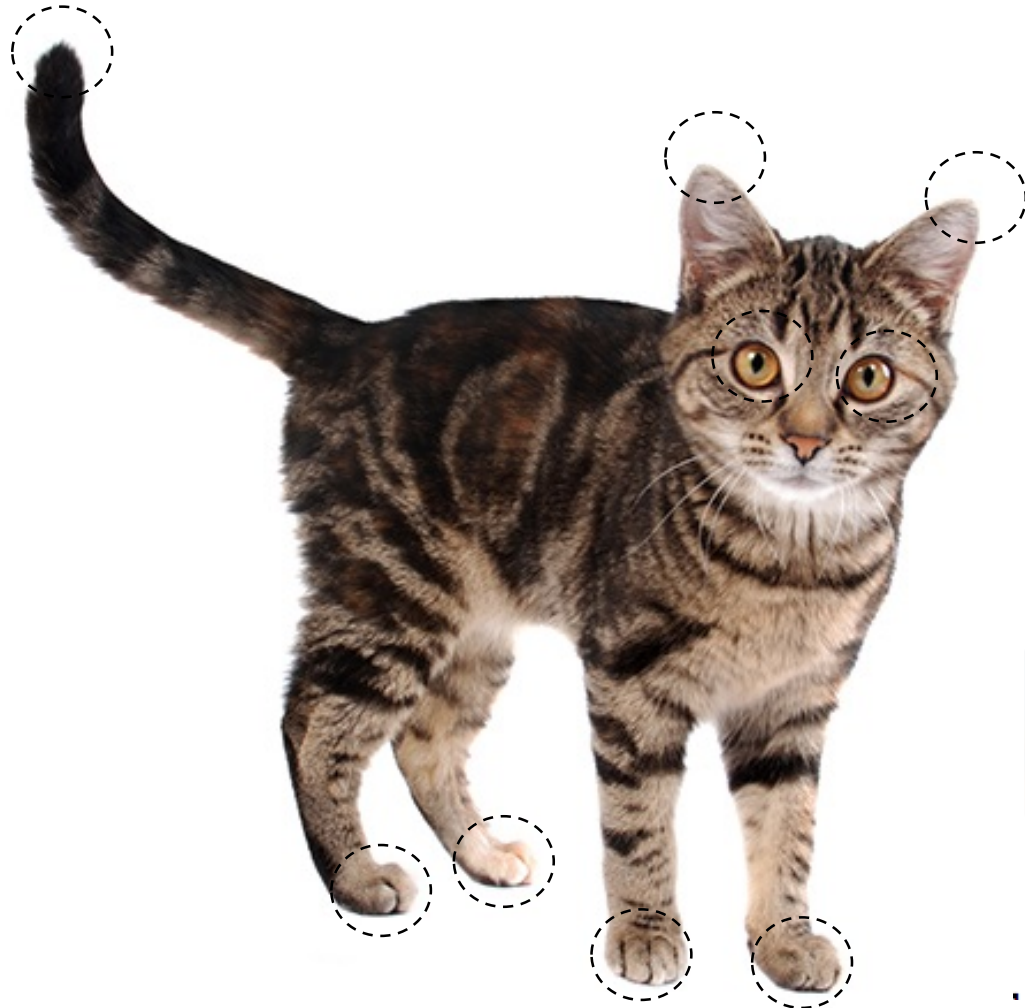




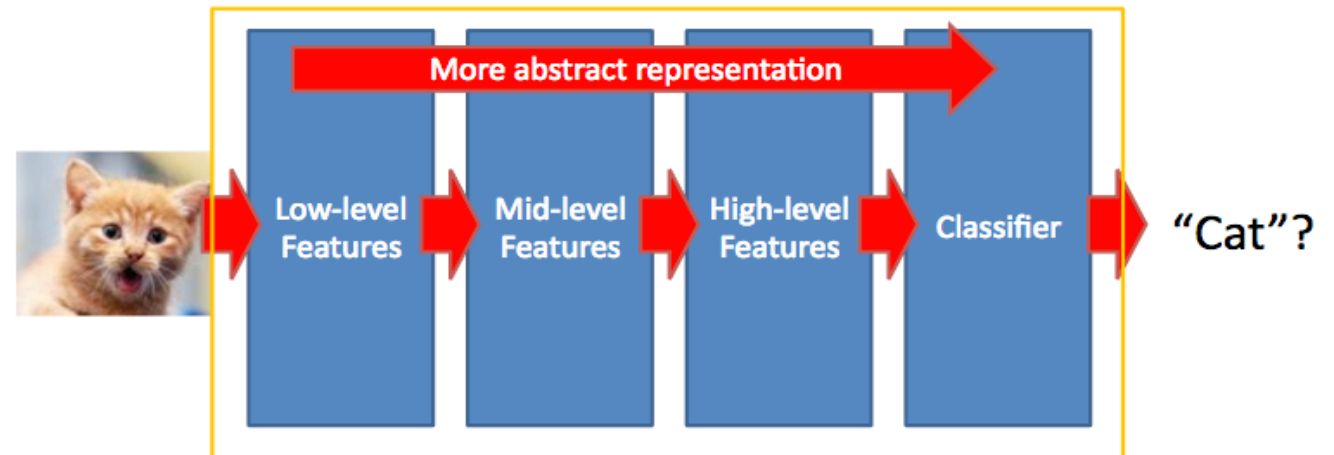
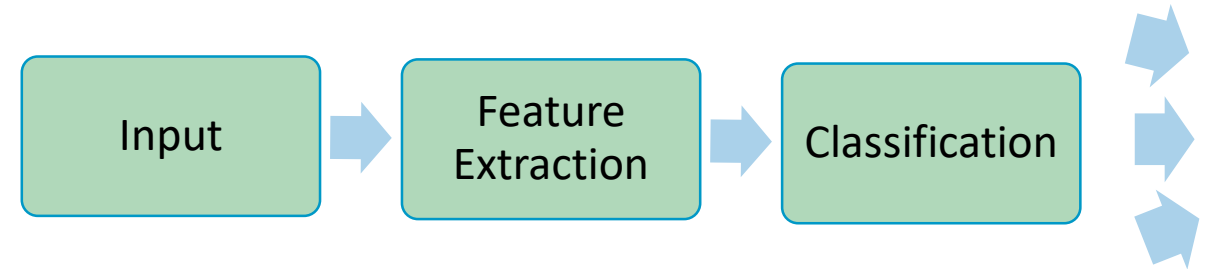
Classification refers to identifying the category to which a new observation belong based on previous examples



Classification refers to identifying the category to which a new observation belong based on previous examples

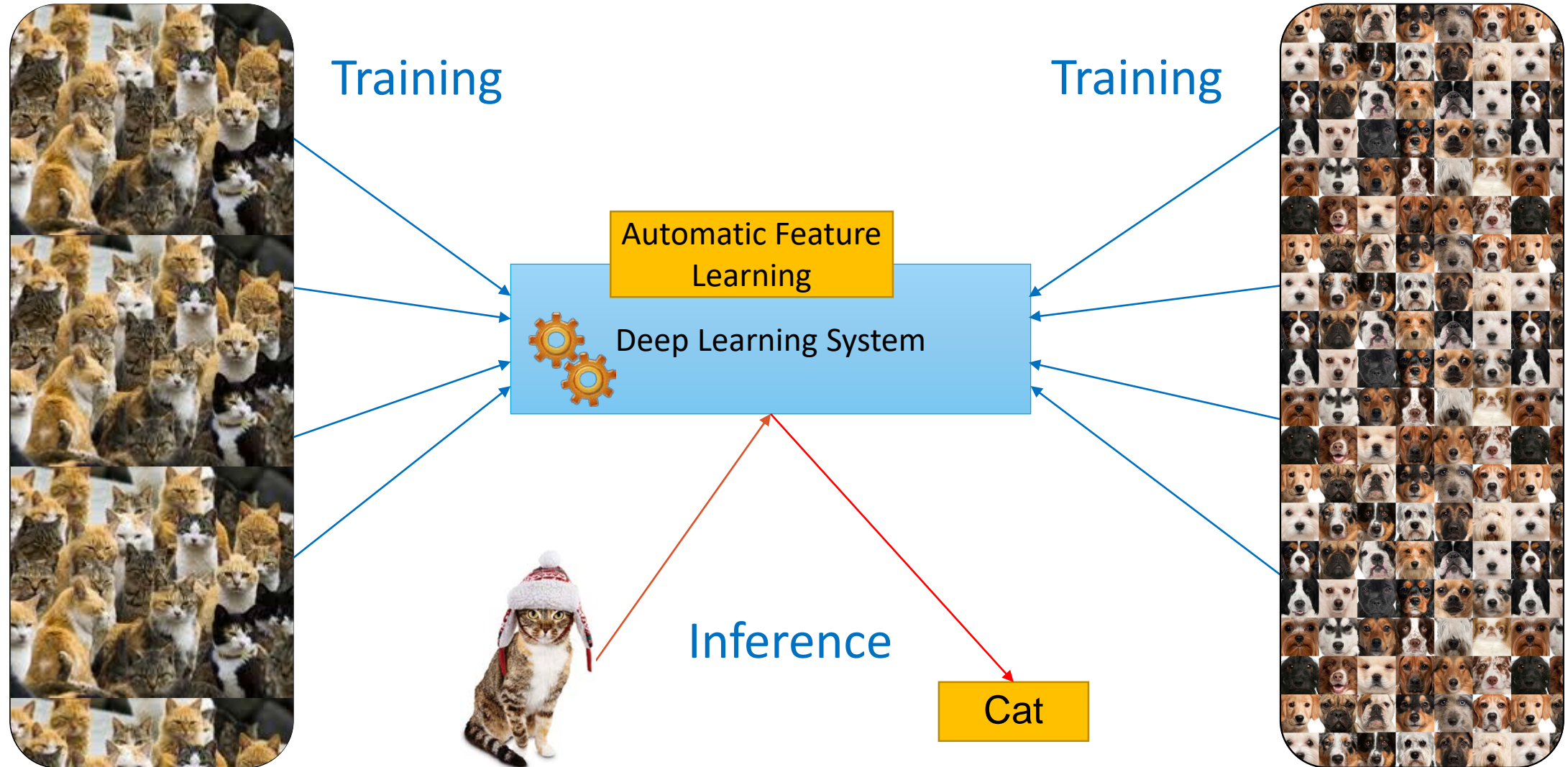


### Typical ML Flow



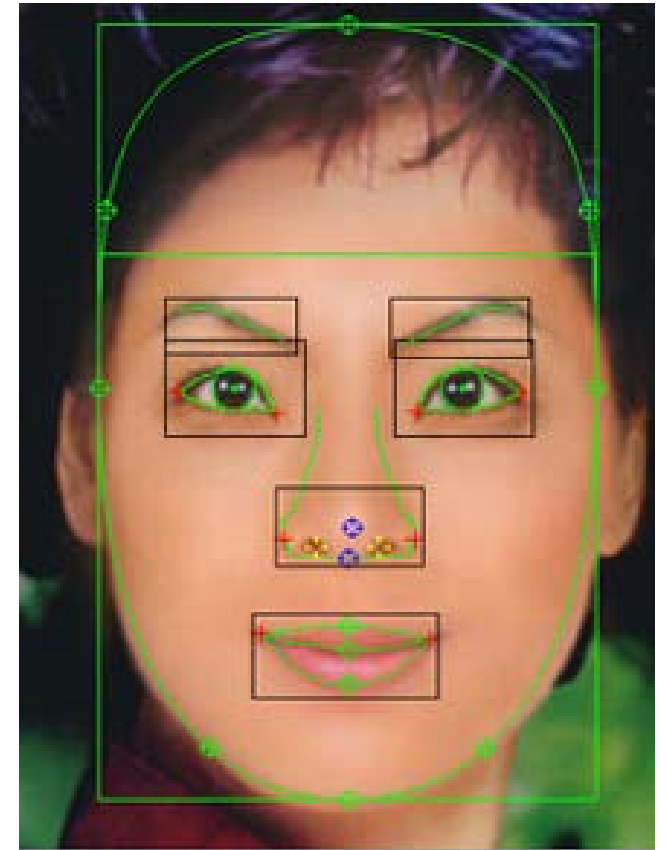
Deep Learning: train layers of features so that classifier works well.

Deep learning is a machine learning technique in which identification and extraction of features is directly done from the dataset





- Deep learning to create automatically features from data
- Features describe the data content in an abstract representation
- For example face features, edges, color distribution, time series, etc.



# Three pillars of deep learning

- Algorithms
- Data
- Computational power

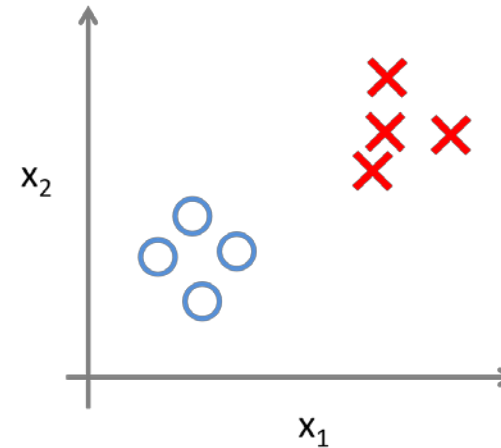


# Supervised and unsupervised learning

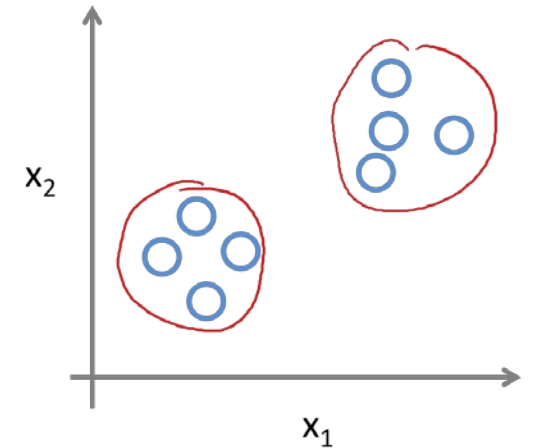
- Supervised

- Labelled training data
- Algorithm learns from data

Supervised Learning



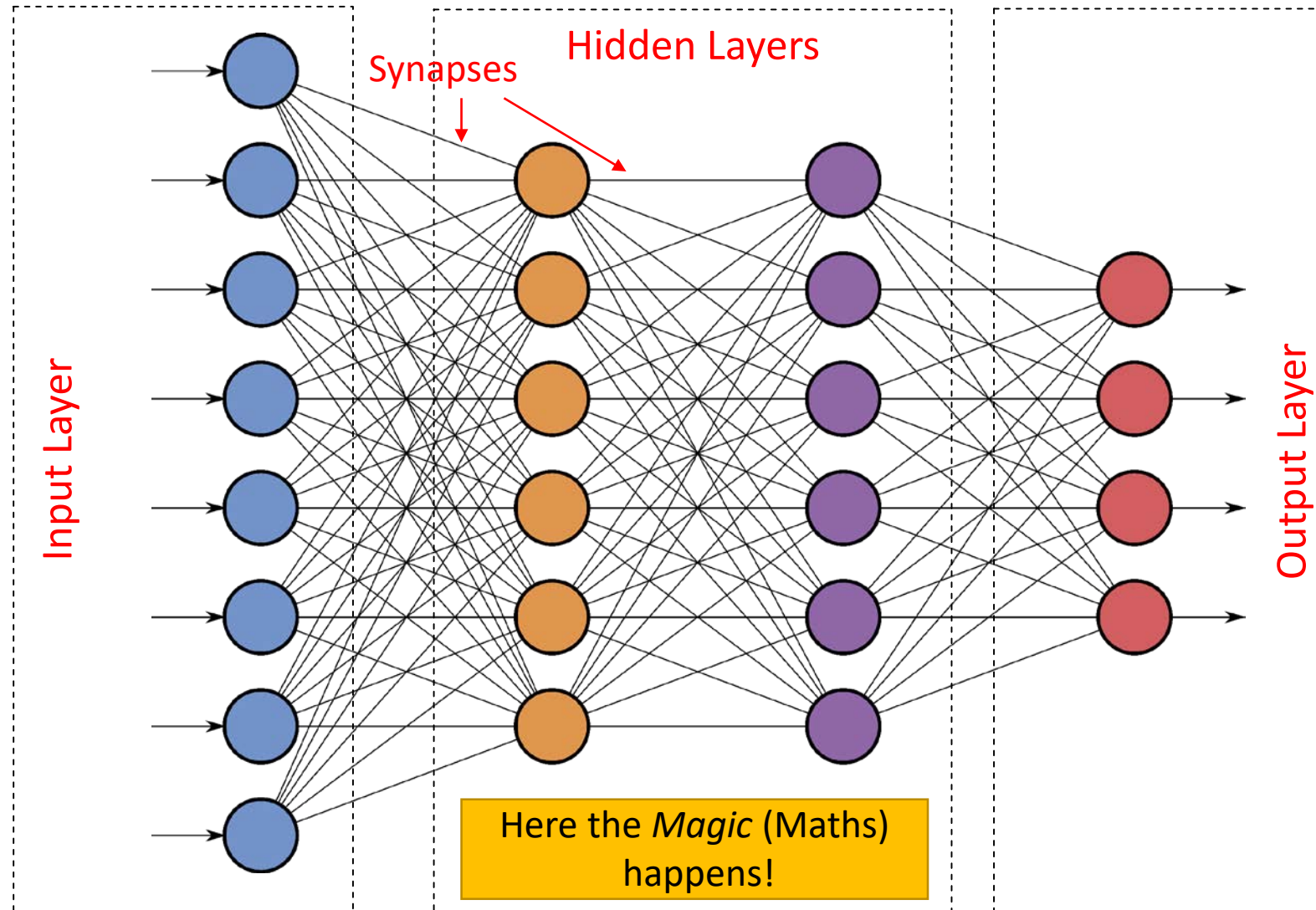
Unsupervised Learning



- Unsupervised

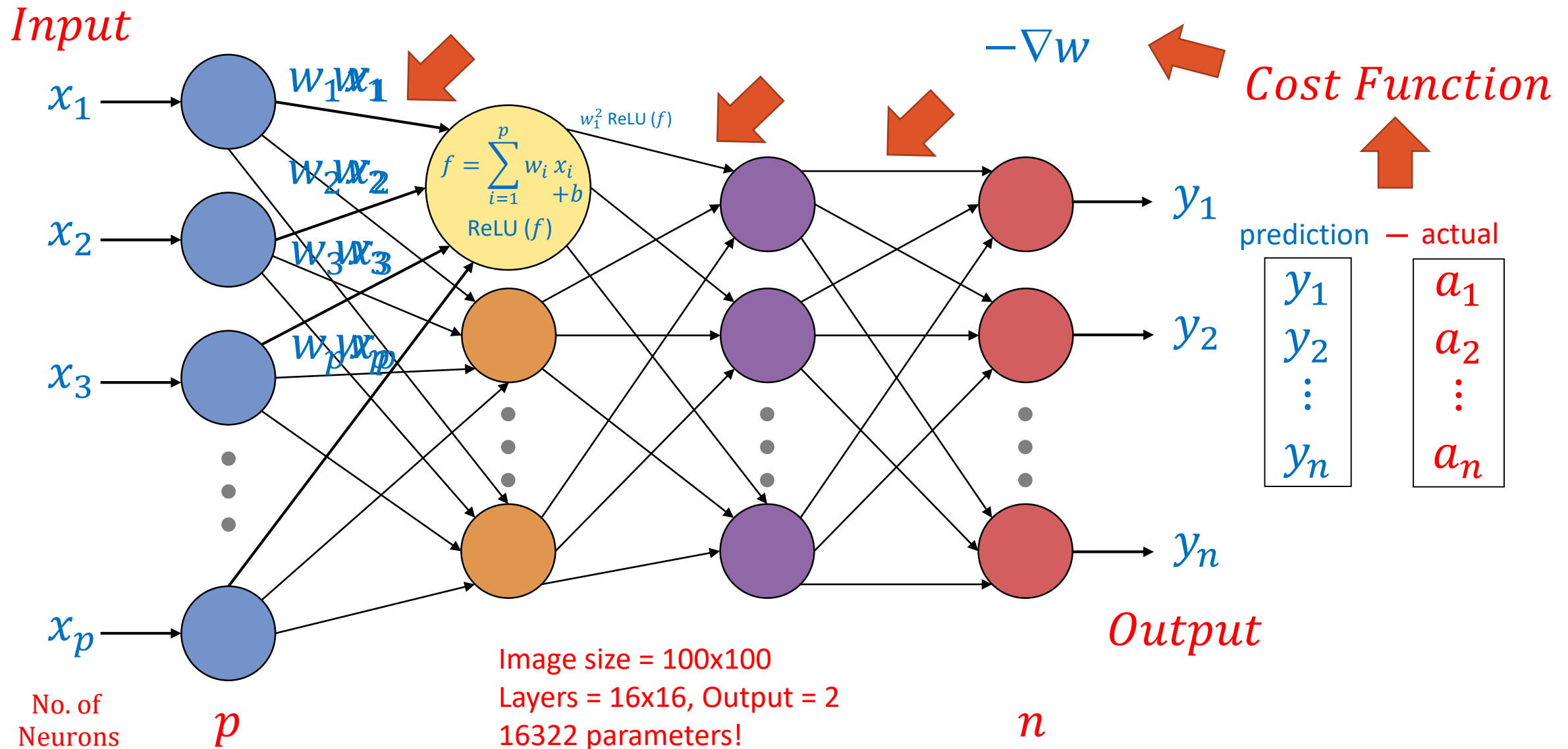
- No training data
- Algorithm learns by itself

Deep learning is predominately based on artificial neural networks (ANNs).  
ANNs progressively *learn* using *neurons* arranged in many layers



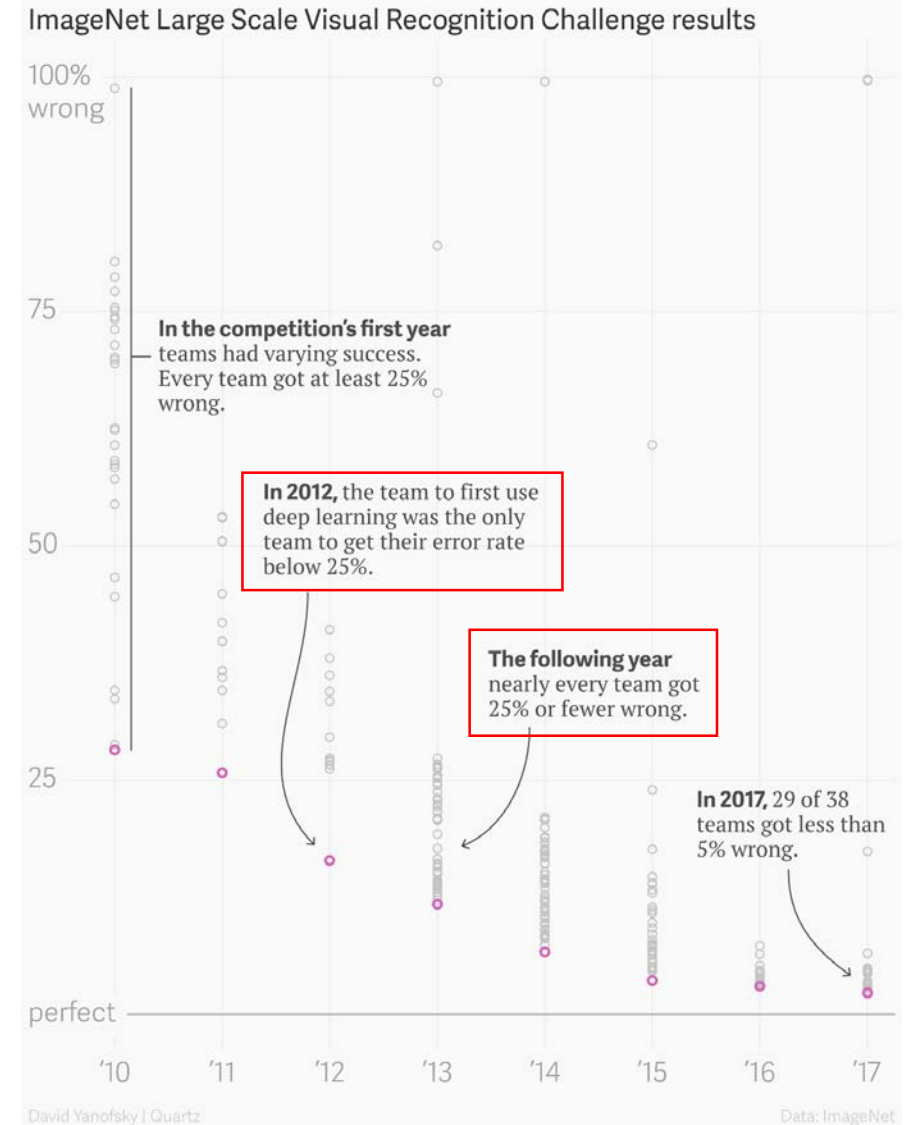
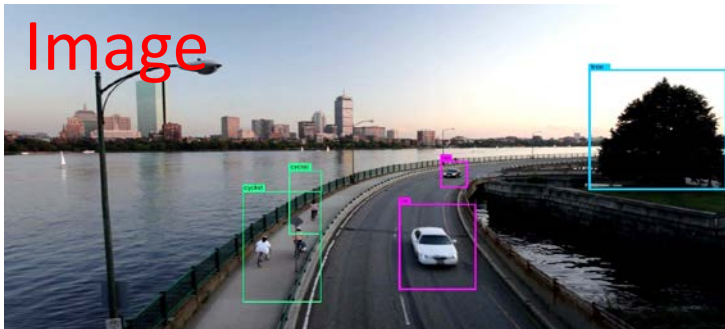


In training neural networks, prediction errors are *backpropagated* for gradually tuning weights and biases to values that predict better

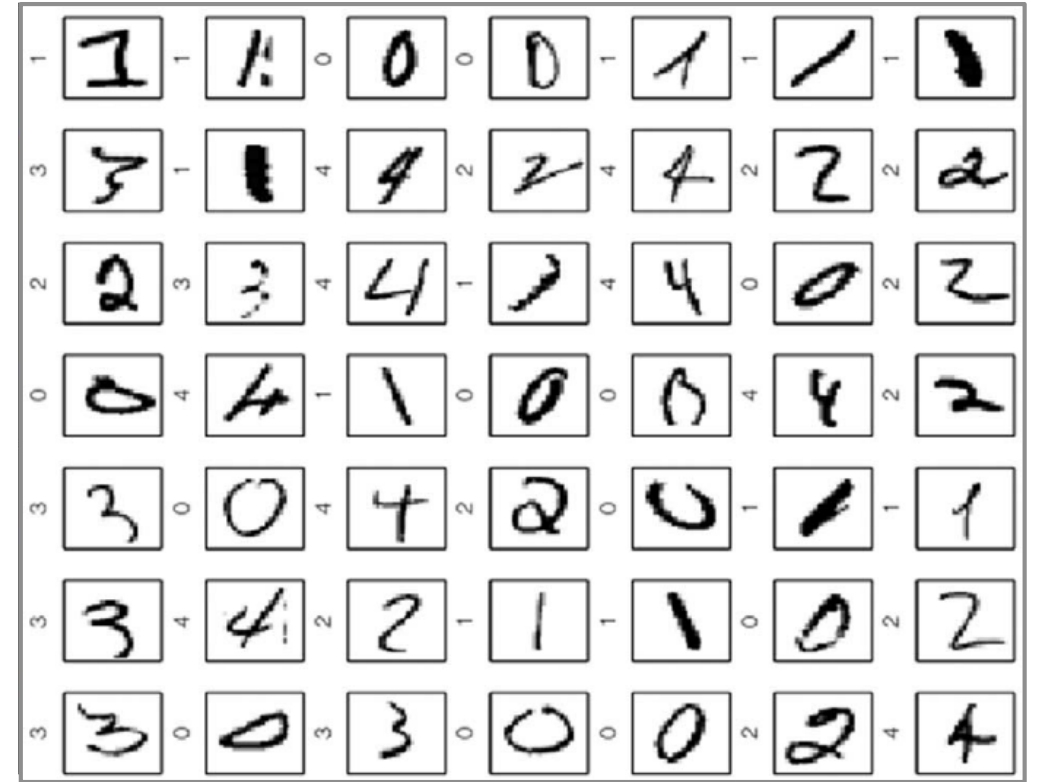
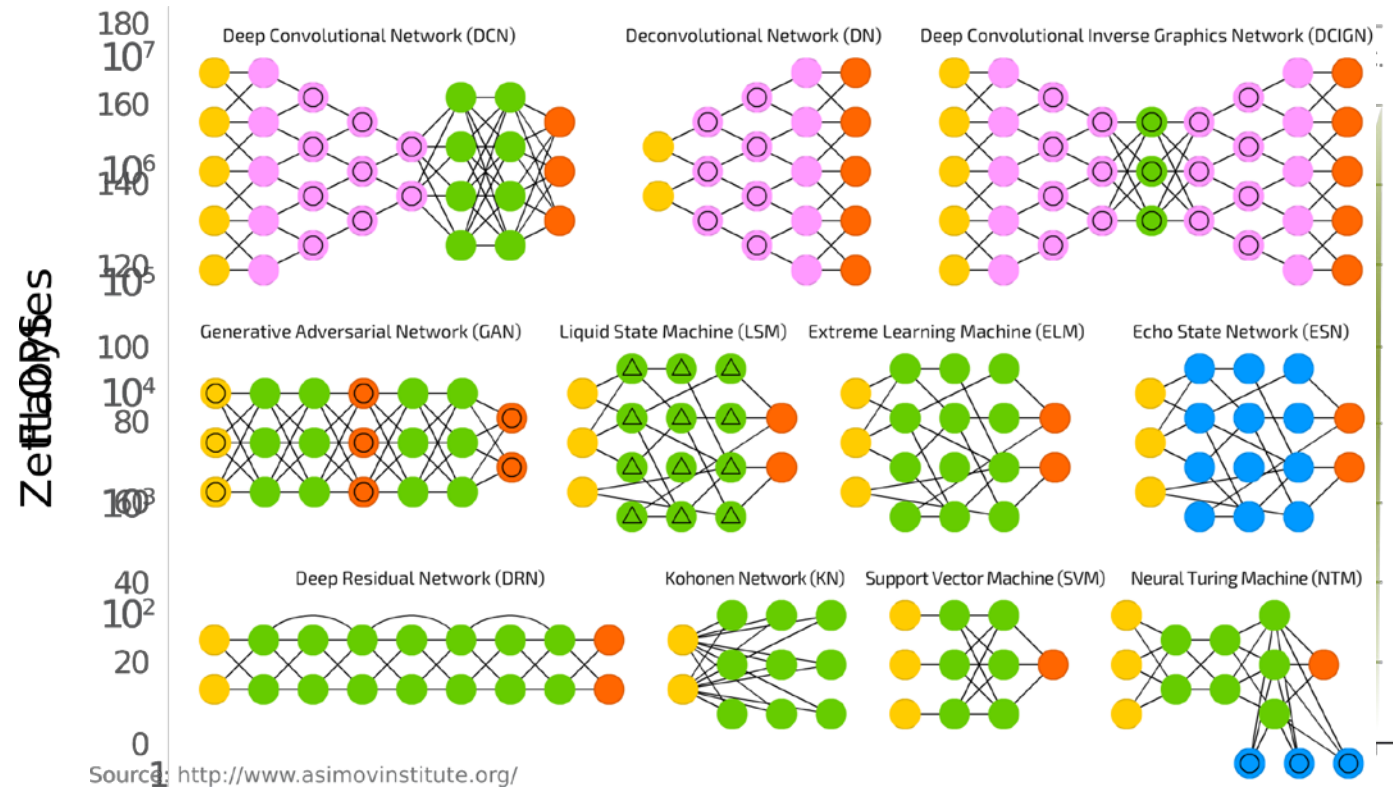


# Deep learning has greatly improved prediction accuracies in various fields, and has become a backbone of artificial intelligence

- Improving prediction accuracies
  - Important in real-world applications such as disease identification, self-driving cars
- Fueling new applications
  - Science, Finance, Healthcare, Automotive



# Availability of large amount of data, improved processing capabilities, and advances in neural networks contributed to the deep learning success



ImageNet is a large database of (URLs) of about 14 million images with more than 20000 different classes.

- Availability of internet and data capabilities made it possible to execute
  - Digital universal networks, Recurrent Neural Networks, Long Term-short memory (LSTM)
- Availability of large datasets
  - ResNet, VGG19, InceptionV3, InceptionResNetV2
  - Tesla V100 capable of 15 TFLOPS of single-precision operations

# Large, complex models and training using large datasets improves classification accuracy of deep learning algorithms

- State-of-the-art accuracies are achieved using very large models
  - Millions / Billions of parameters, Large number of layers and depth
    - Neural network of Google Brain<sup>1</sup> has 137 billion parameters!

Size and depth of the model, Top-1 accuracy, ImageNet dataset<sup>2</sup>

Model	Top-1 Accuracy %	No. of Parameters	Network Depth
MobileNet	66.5	4,253,964	88
ResNet50	75.9	25,636,712	168
Xception	79.0	22,910,480	126
VGG16	71.5	138,357,544	23
VGG19	72.7	143,667,240	26
InceptionV3	78.8	23,851,784	159
InceptionResNetV2	80.4	55,873,736	572

Note: Recently, Google's NASNet<sup>3</sup> reached accuracy of 82.7% with 22.6 Million parameters.

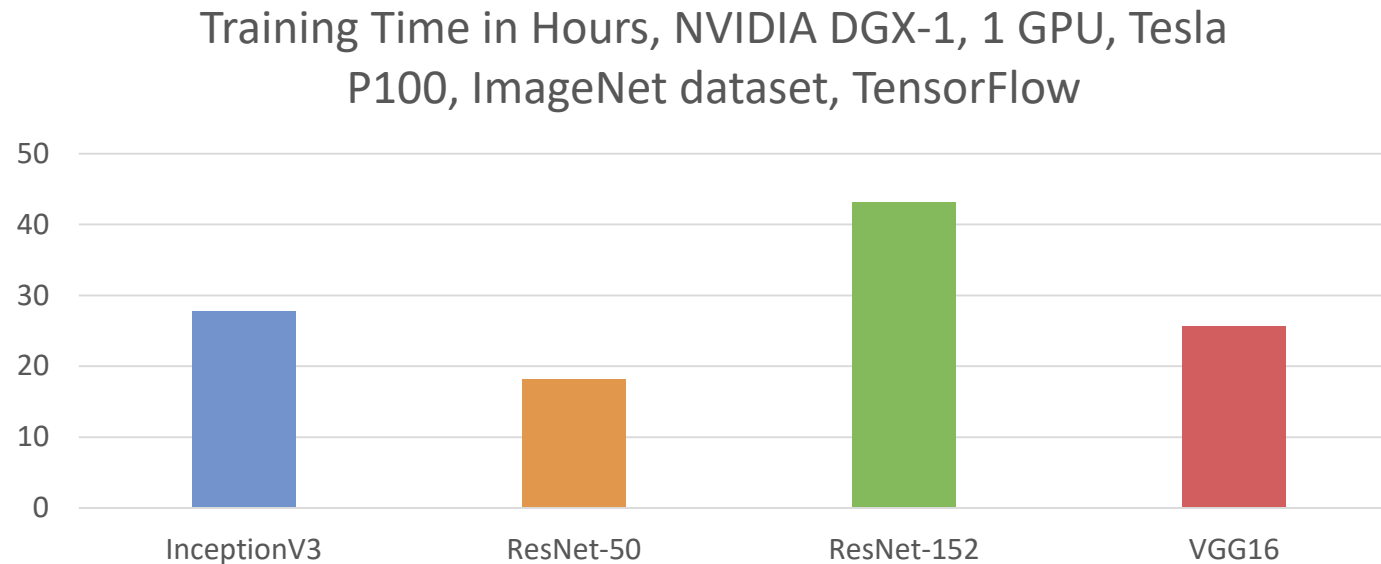
<sup>1</sup>Shazeer, Noam, et al. "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer." arXiv preprint arXiv:1701.06538, 2017.

<sup>2</sup><https://keras.io/applications/> <sup>3</sup>Zoph, Barret, et al. "Learning transferable architectures for scalable image recognition." arXiv preprint arXiv:1707.07012, 2017.



# Efficient distributed deep learning across a large number of nodes is critical for modern deep learning applications

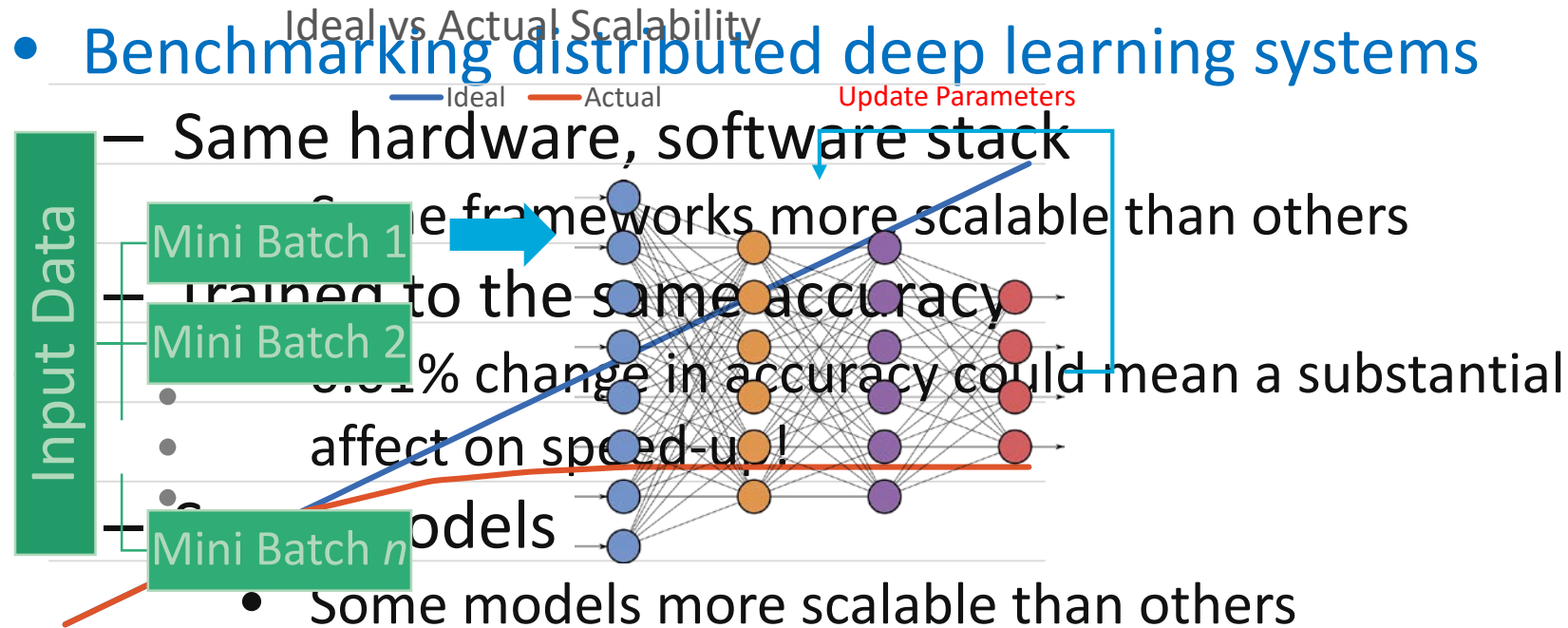
- It takes hours, and even days to train today's models on a single node
  - The models will only become more bigger and complex
    - The amount of training data will only increase
  - Increase in GPU processing power cannot cope up
    - Billions of parameters will not fit in GPU memory



Data: <https://www.tensorflow.org/performance/benchmarks>

# Distributed deep learning should be scalable, should not affect the accuracy of the classification, and should fully utilize available resources

## • Benchmarking distributed deep learning systems



— Due to computation / communication ratio

## • Scaling Efficiency and Communication Overhead

- Ratio between training time of one iteration on 1 nodes to total training time when distributed over  $n$  nodes
  - A larger batch size is usually proved beneficial for distributed deep learning
- Number of epochs required for convergence
  - Reporting affected by  $n$  nodes with 1x GPU vs  $\frac{n}{8}$  nodes with 8x GPUs
- Learning rate adaptation, trade-off between runtime and accuracy

Goyal, Priya, et al. "Accurate, Large Mini-batch SGD: Training ImageNet in 1 Hour." arXiv preprint arXiv:1706.02677, 2017.

**This presentation will introduce distributed deep learning, walk through prominent techniques, and identify existing challenges and future directions**



## **Introduction and Motivation**

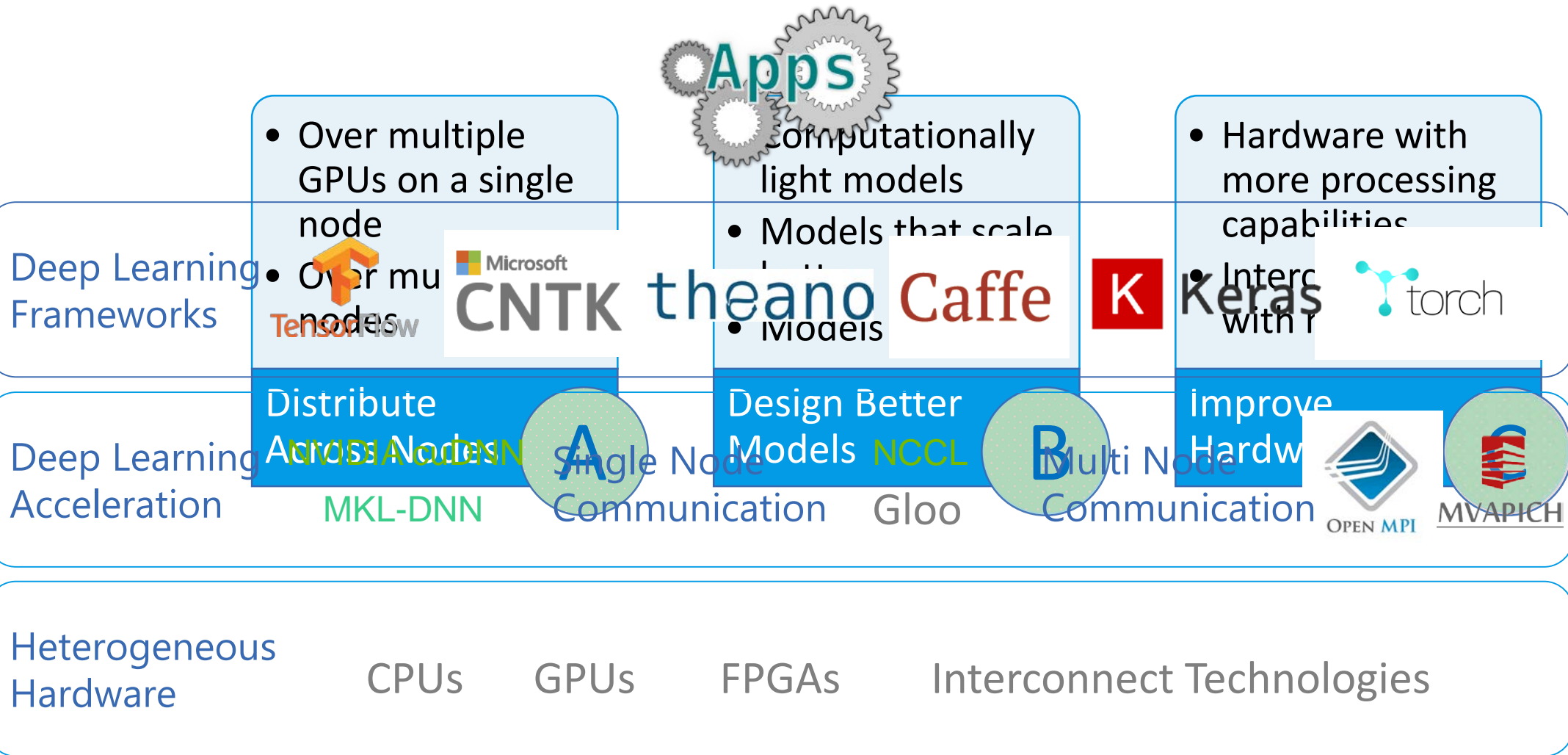


## **Existing Techniques and Toolsets**



## **Future Directions**

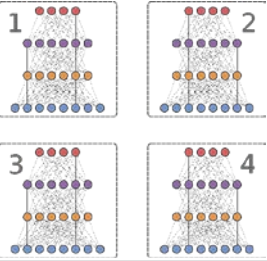
# Scaling efficiency of distributed deep learning can be improved at different levels, and a co-design approach is needed





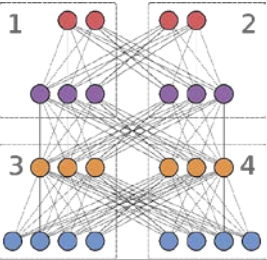
# Distribute Across Nodes

# There are three parallelization methods employed in distributed implementation of deep learning applications



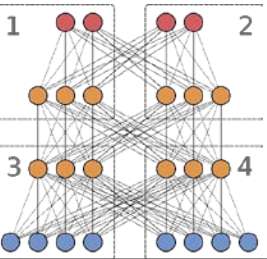
## Data Parallelism

- Data partitioned across machines, each machine holds local copy of the model, synchronization required



## Model Parallelism

- Neural network partitioned across machines, Parallelizing mathematical operations across machines



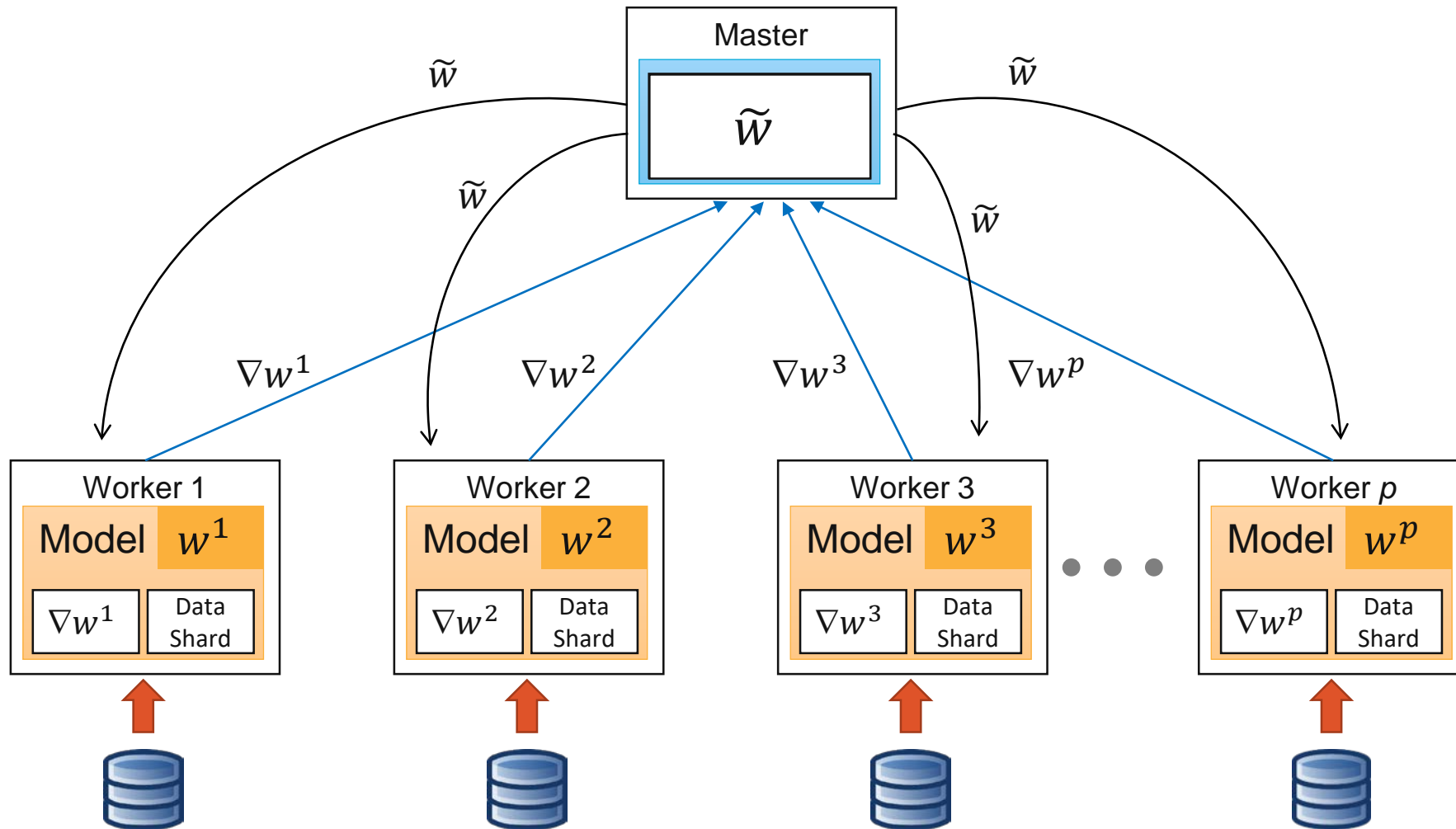
## Hybrid Approaches

- Data partitioning for some parts of neural network, model partitioning for correctness on some parts, Automatic selection

Distribute  
Across Nodes

A

**Data Parallelism: In Synchronized data parallelism, the algorithm includes two parts: sum of local gradients and broadcast of global weight to workers**



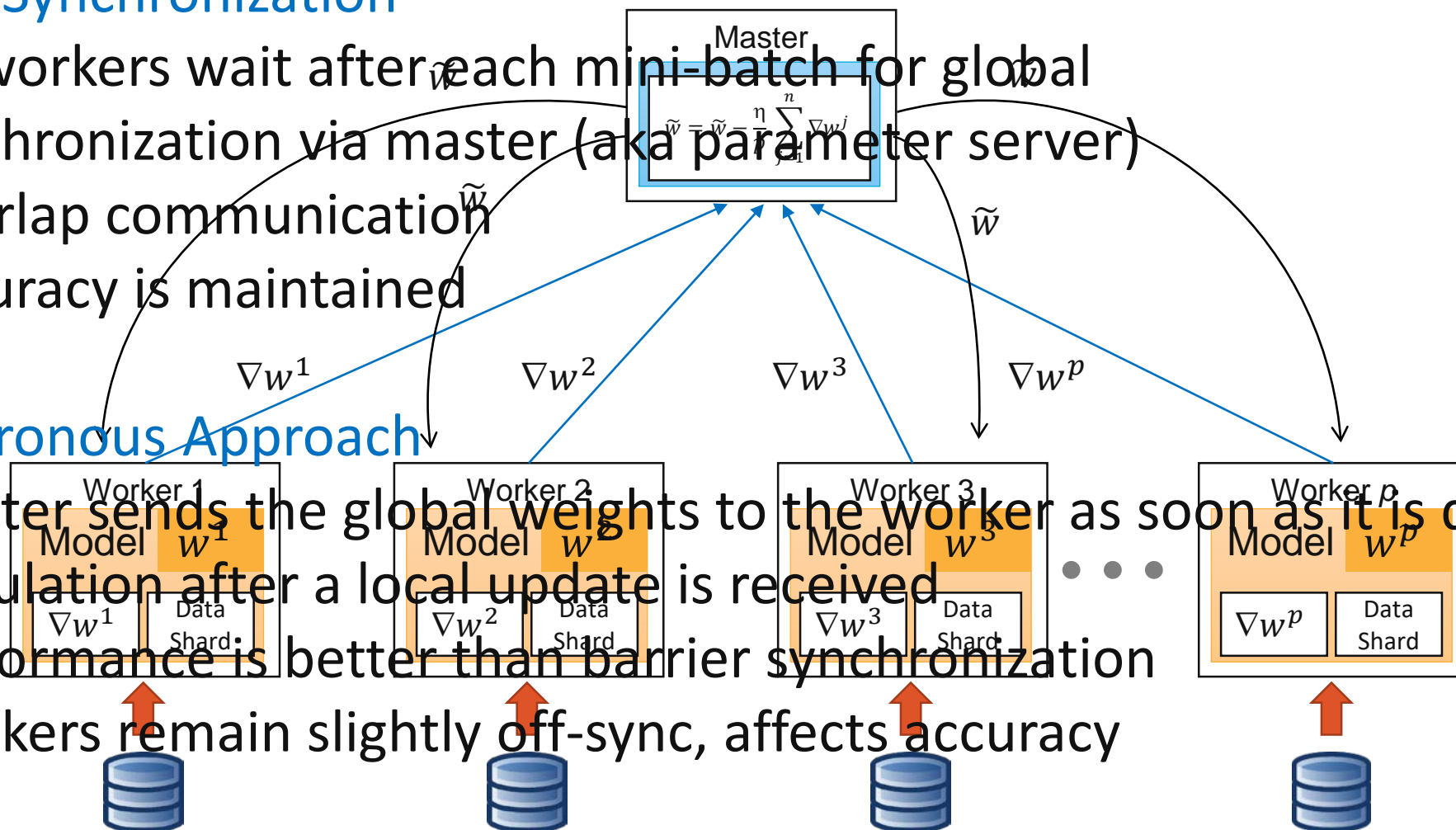
# Data Parallelism: Barrier before each update, overlapping communication and computation, or asynchronous approaches are generally employed

- Barrier Synchronization

- All workers wait after each mini-batch for global synchronization via master (aka parameter server)
- Overlap communication
- Accuracy is maintained

- Asynchronous Approach

- Master sends the global weights to the worker as soon as it is done calculation after a local update is received
- Performance is better than barrier synchronization
- Workers remain slightly off-sync, affects accuracy

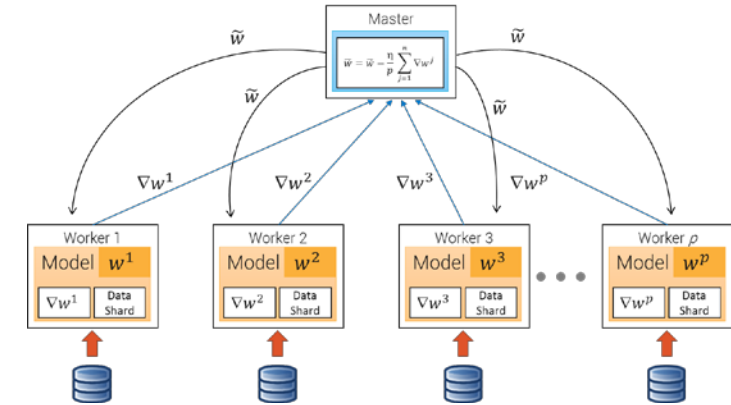
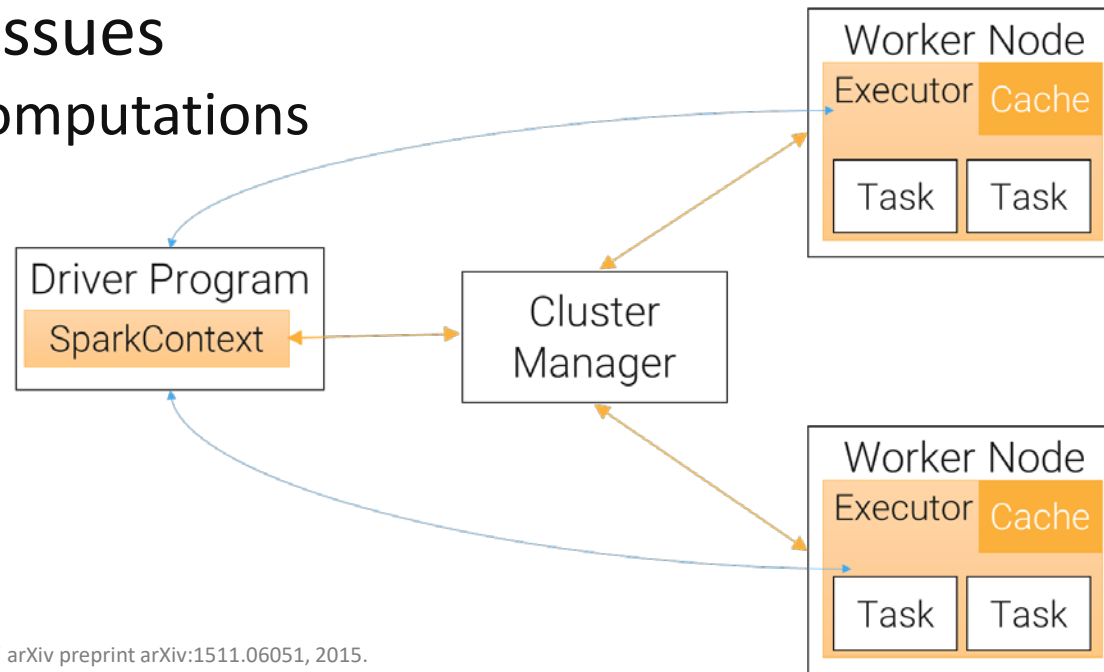




# Data Parallelism: Big data processing frameworks and resource management systems also make good candidates for distributed deep learning

- Spark-based Solutions

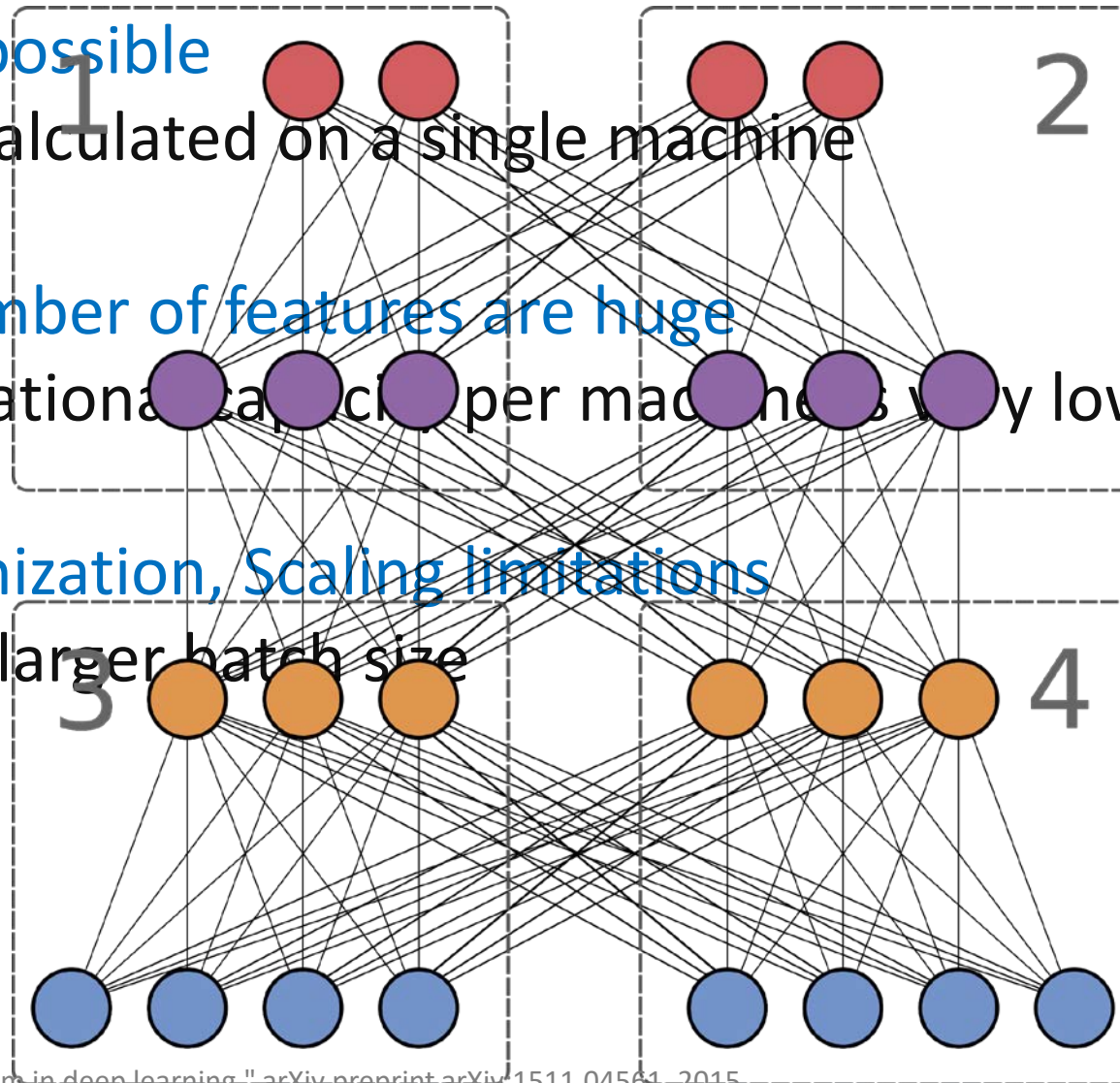
- Integration of deep learning frameworks, e.g. TensorFlow with Spark
  - ClusterManager takes role of parameter server
  - Mesos used as resource manager
- Performance issues
  - Iterative computations



Moritz, Philipp, et al. "Sparknet: Training deep networks in spark." arXiv preprint arXiv:1511.06051, 2015.  
Kim, Hanjoo, et al. "Deepspark: Spark-based deep learning supporting asynchronous updates and caffe compatibility." CoRR, vol. abs/1602.08191, 2016.

# Model Parallelism: Neural network is divided into several mini-models and distributed across multiple nodes

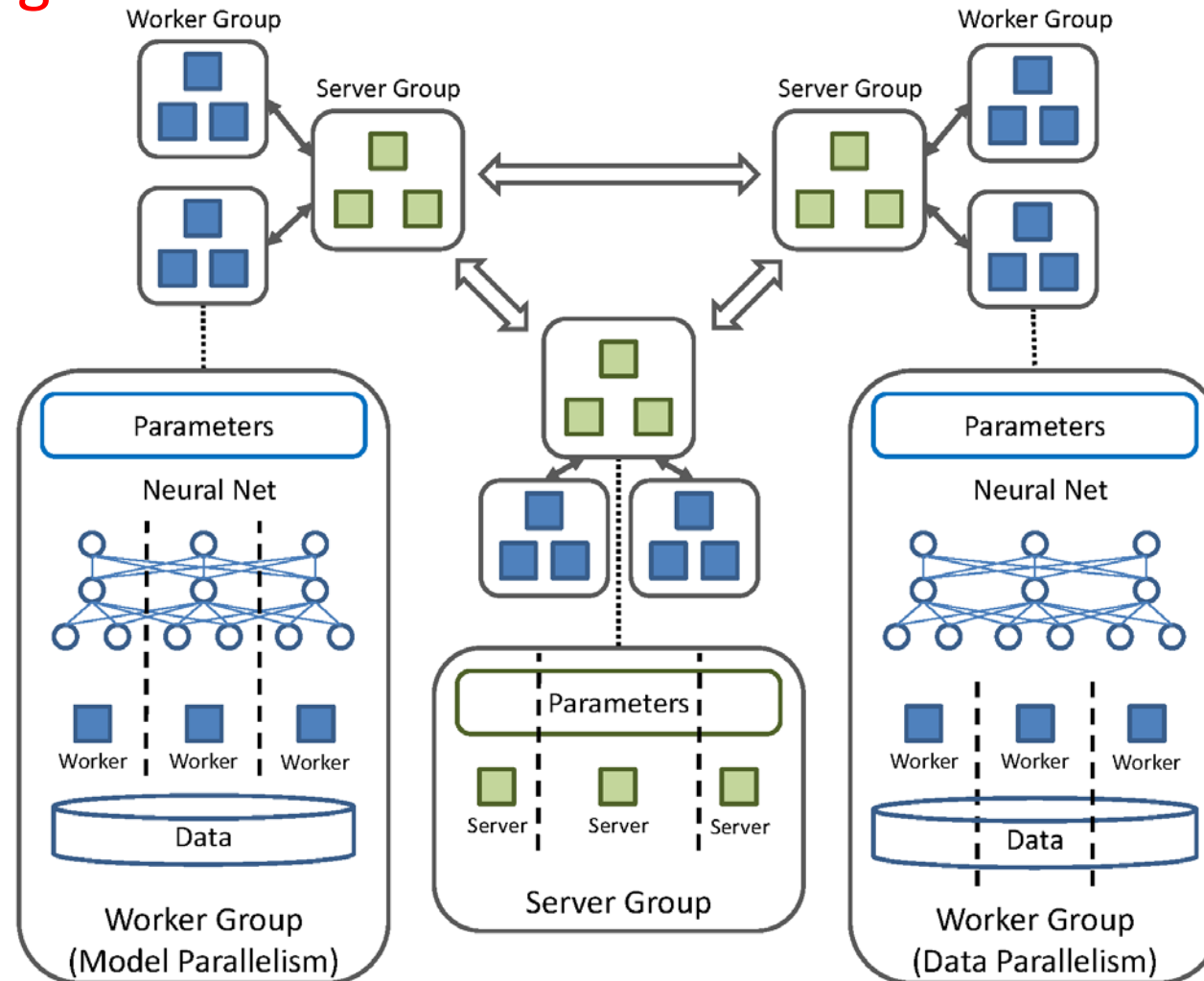
- Highest accuracy possible
  - As if model is calculated on a single machine
- Feasible when number of features are huge
  - Or the computational capacity per machine is very low
- Frequent synchronization, Scaling limitations
  - Problems with larger batch size



Dettmers, Tim. "8-bit approximations for parallelism in deep learning." arXiv preprint arXiv:1511.04561, 2015.

# Hybrid Approach: Some parts of the neural network are divided using model parallelism while other parts employ data parallelism

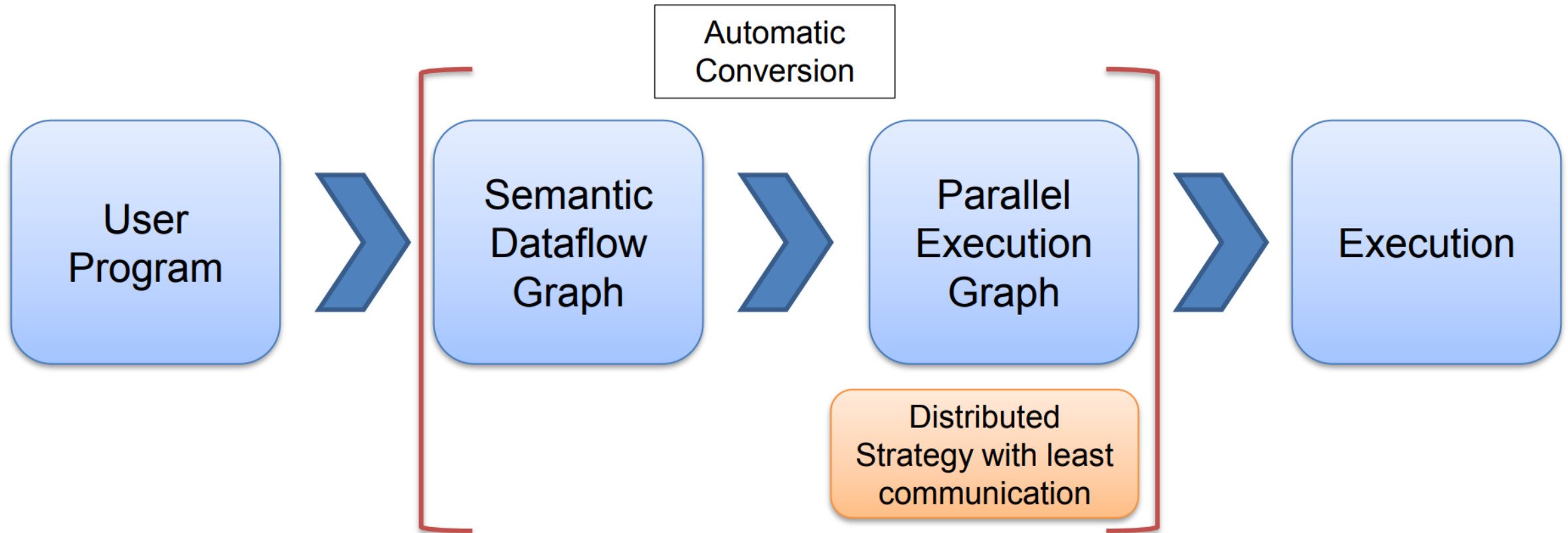
## Example: Apache Singa



Source: <http://singa.incubator.apache.org/>

# Hybrid Approach: Automatic selection of the best possible parallelism model has also been proposed in literature

## Example: Tofu



Automatic generation of Neural networks using evolutionary algorithm, ORNL.

Source: Tofu – Parallelizing Deep Learning Systems with Automatic Tiling, Minjie Wang, 2017

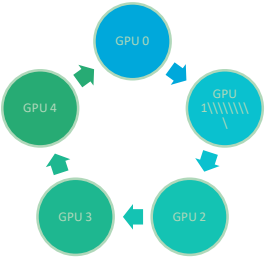
# Communication efficiency plays an important role in the scalability of the distributed deep learning algorithms

Communication overhead is the difference between runtime of one iteration when distributed over  $n$  processing units and runtime of one iteration on a single unit



## Interconnect Bandwidth and Latency

- Inter-node communication, Intra-node communication

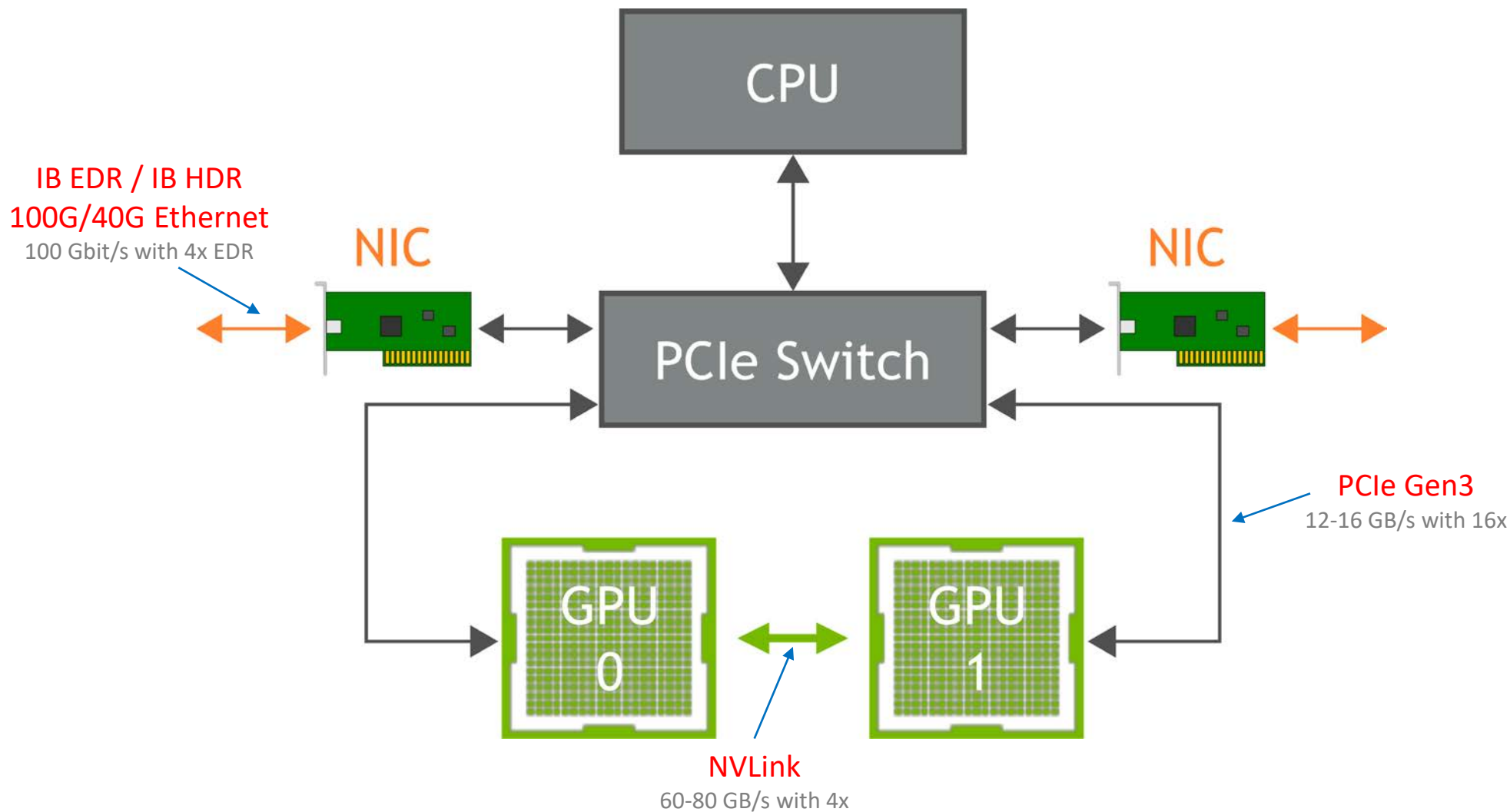


## Communication Algorithm and Libraries

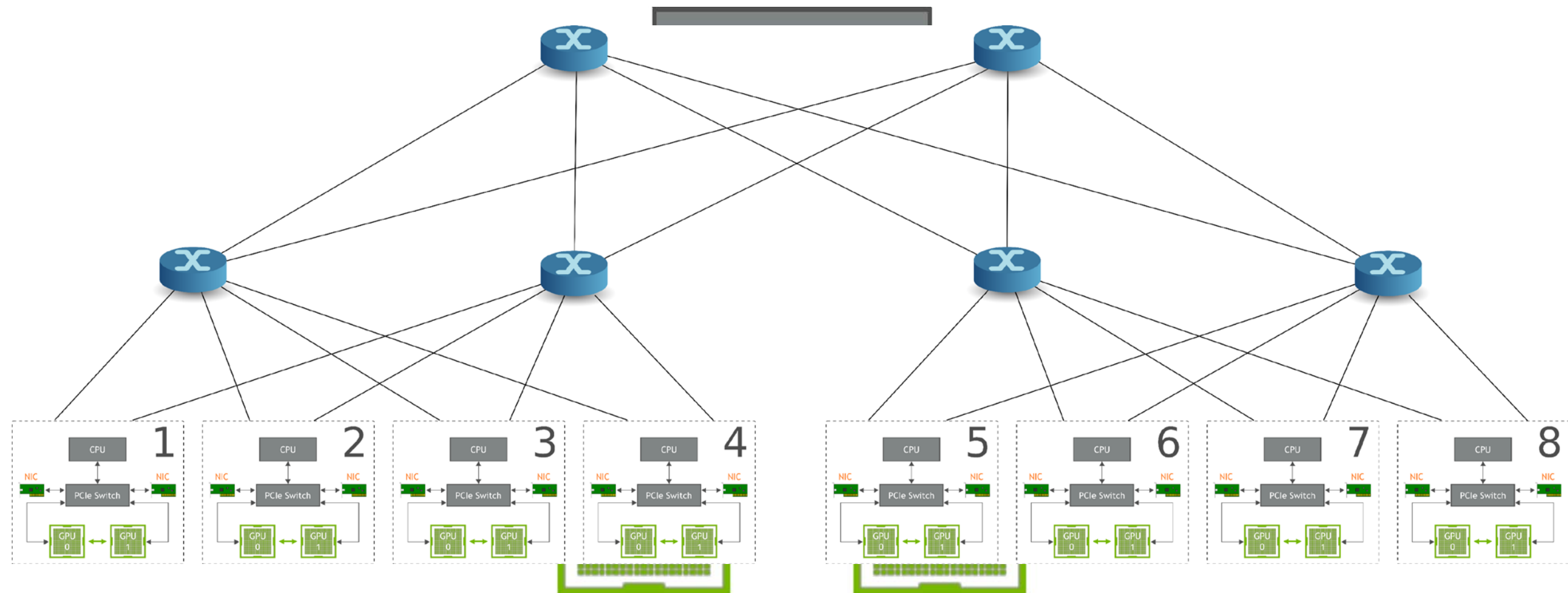
- Topology-aware, Adaptive to the configuration



In large-scale systems, a hierarchy of interconnect technologies are used giving different bandwidths, latencies, and connectivity



# In large-scale systems, a hierarchy of interconnect technologies are used yielding different bandwidths, latencies, and connectivity



- In such distributed scenarios, topology and routing also plays an important role
  - Communication algorithm must adapt

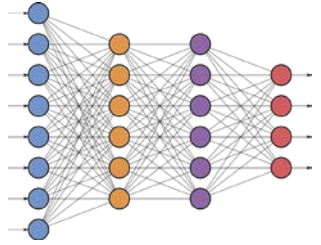
MLLink  
60-80 GB/s with 4x

# Communication algorithm can optimize utilization of the available link capacity by avoiding contention in data transfers

- Collective operations
  - All-Reduce operations in synchronous data parallelism
  - All-Gather and broadcast in asynchronous data parallelism / model parallelism
- Communication Libraries
  - MPI is very mature, standardized, and provide excellent scale-out performance
    - Scale-Up?
  - NICCL
    - Scale Up, Across multiple nodes is an issue even though supported now

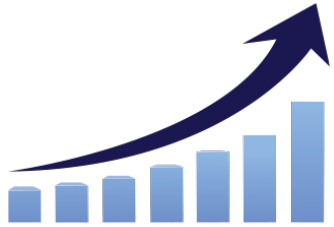
# Design Scalable Models

# Distributed deep learning can be accelerated through the use of computationally lighter & scalable models and define-by-run methodologies



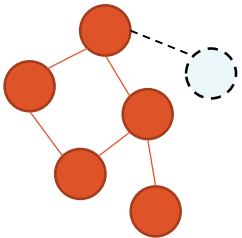
## Computationally Lighter Models

- Some models are inherently more computationally expensive than others, VGG vs Inception



## Scalable Models

- Models that are more scalable, tolerant to delayed parameter updates



## Define-by-Run Paradigm

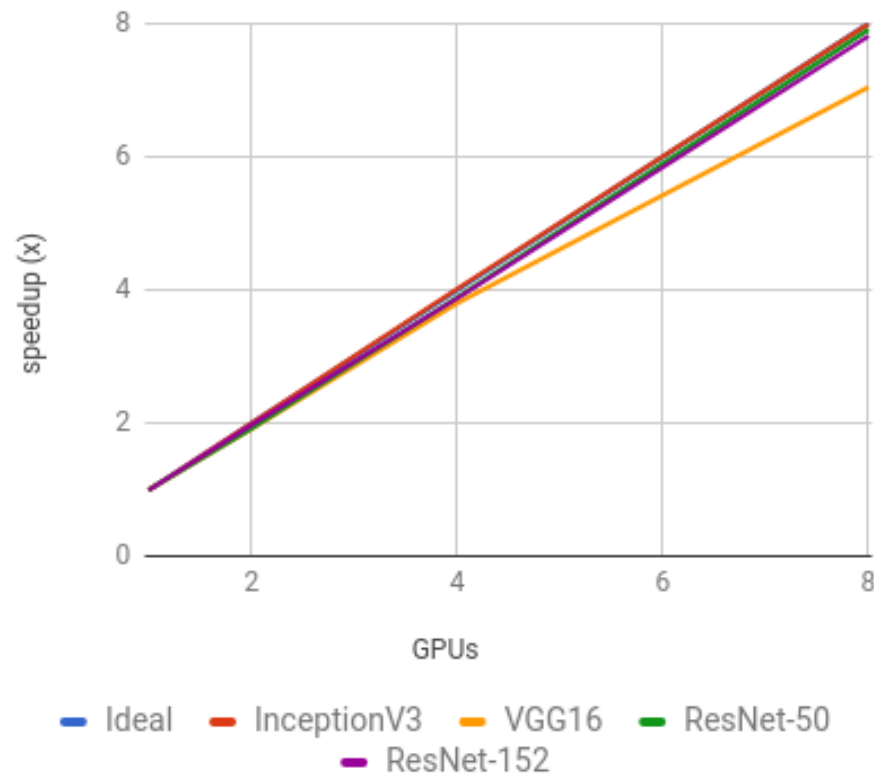
- Dynamic graph updates, Potential for adaptation to the changing data



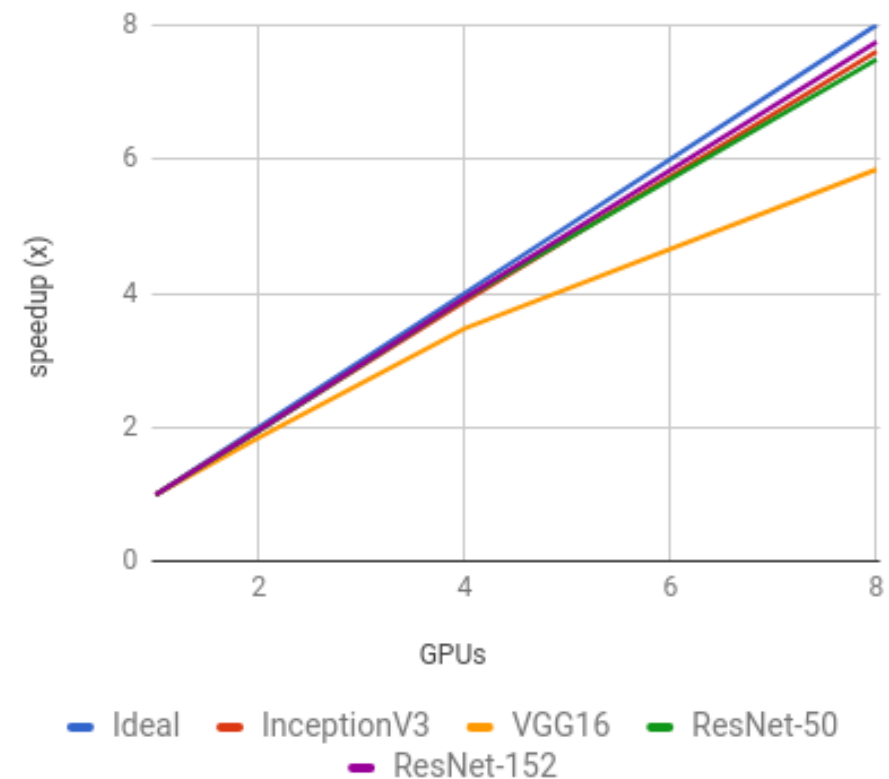
# Models can also be designed in a way that they are more scalable and suited for the distributed deep learning

## Example: VGG16 vs ResNet-152 – Single node, ImageNet, TensorFlow

Tesla® P100 speedup (synthetic data)



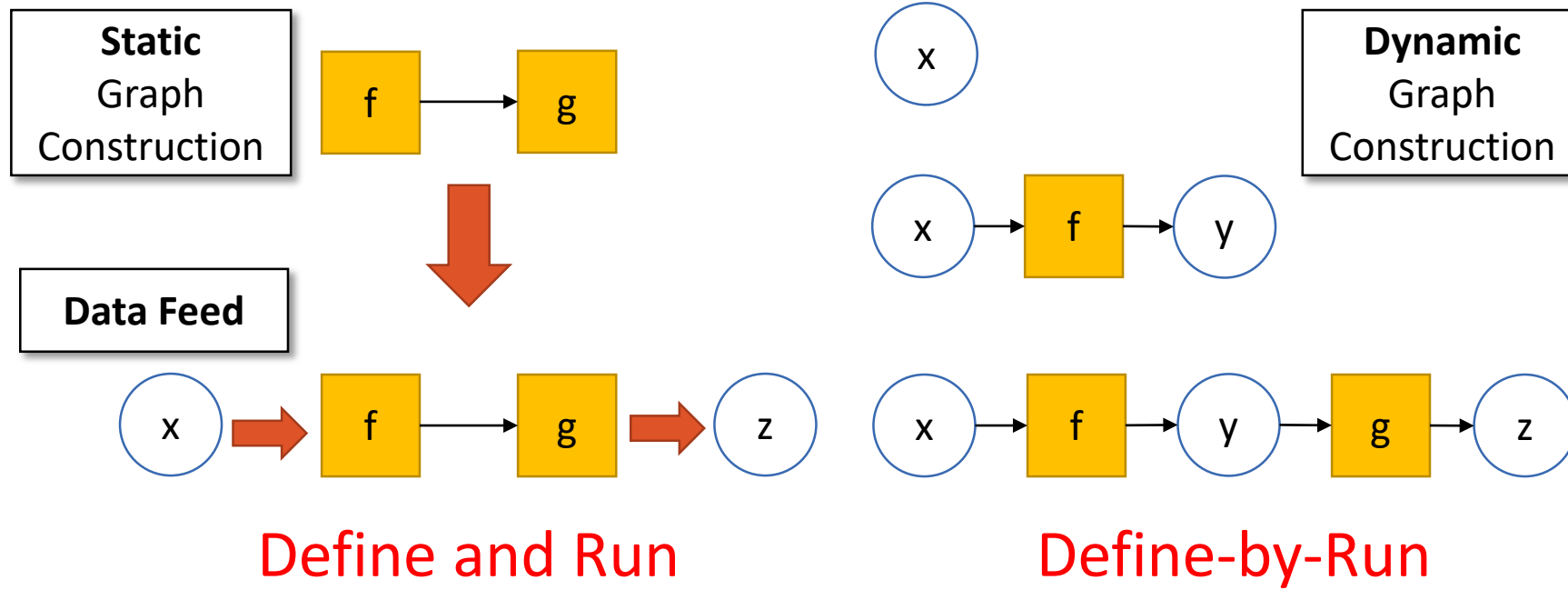
Tesla® P100 speedup (real data)



Source: <https://www.tensorflow.org/performance/benchmarks>

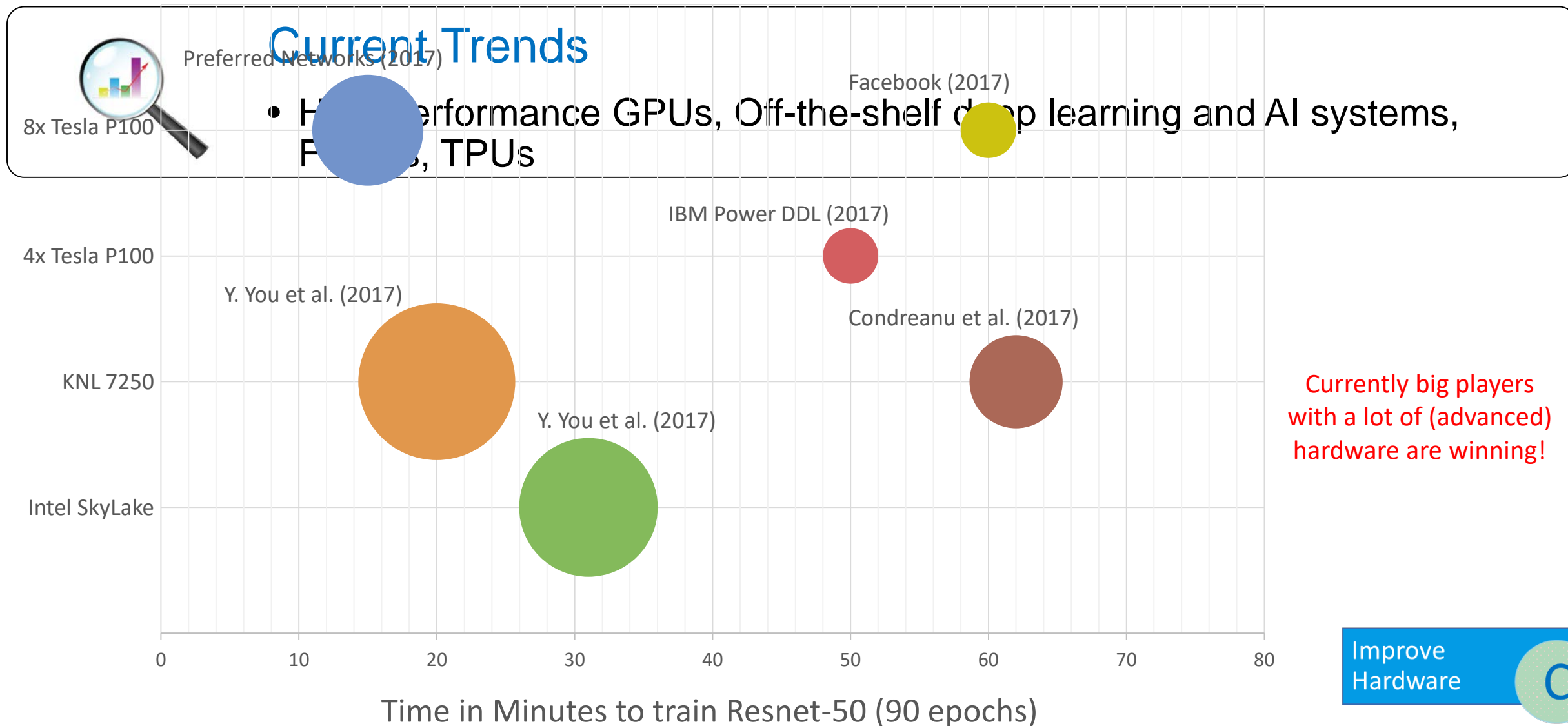
# In Define-by-Run frameworks, graphs constructions *on the fly* giving more flexibility and potential performance improvements

- Chainer
  - The first define-by-run deep learning framework
- TensorFlow
  - Eager execution



# Improve Hardware

# Hardware advances will still play an integral role in improving efficiency of deep learning algorithms



This presentation will introduce distributed deep learning, walk through prominent techniques, and identify existing challenges and future directions



## Introduction and Motivation



## Existing Techniques and Toolsets



## Future Directions

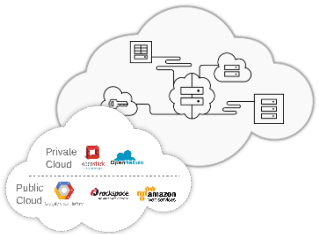


# Future research directions include a holistic approach for variety of workloads & heterogeneous environments, and NN-aware communication



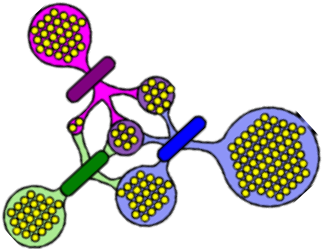
## Holistic Approach

- New workloads are emerging, different toolsets, in modern apps workload of different types will co-exist



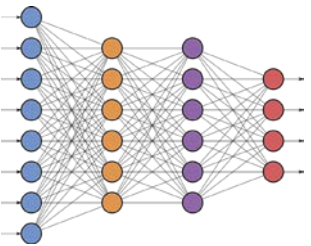
## Heterogeneous Environments

- Accelerators, Adaptive to the configuration, Dynamic resource allocations, Clouds (challenges related to multi-tenancy)



## Neural Network Aware Communication

- Pipelining communication with computation, Topology-aware, Offloading



## Improved Algorithms and Models

- Less communication requirements, Good accuracy with large batch sizes, adaptive learning rates, Automatic neural network

# References – Incomplete list

Dean, Jeffrey, et al. "Large scale distributed deep networks." Advances in neural information processing systems. 2012.

Cho, Minsik, et al. "PowerAI DDL." arXiv preprint arXiv:1708.02188 (2017).

Goyal, Priya, et al. "Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour." arXiv preprint arXiv:1706.02677 (2017).

<https://www.tensorflow.org/>

Akiba, Takuya, Shuji Suzuki, and Keisuke Fukuda. "Extremely Large Minibatch SGD: Training ResNet-50 on ImageNet in 15 Minutes." arXiv preprint arXiv:1711.04325 (2017).

Shazeer, Noam, et al. "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer." arXiv preprint arXiv:1701.06538 (2017).

Chilimbi, Trishul M., et al. "Project Adam: Building an Efficient and Scalable Deep Learning Training System." OSDI. Vol. 14. 2014.

Abadi, Martín, et al. "Tensorflow: Large-scale machine learning on heterogeneous distributed systems." arXiv preprint arXiv:1603.04467 (2016).

Tokui, Seiya, et al. "Chainer: a next-generation open source framework for deep learning." Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference on neural information processing systems (NIPS). Vol. 5. 2015.

Jozefowicz, Rafal, et al. "Exploring the limits of language modeling." arXiv preprint arXiv:1602.02410 (2016).

<https://www.nvidia.com/en-us/deep-learning-ai/>

<https://keras.io/>

Bottou, Léon. "Large-scale machine learning with stochastic gradient descent." Proceedings of COMPSTAT'2010. Physica-Verlag HD, 2010. 177-186.

Zoph, Barret, et al. "Learning transferable architectures for scalable image recognition." arXiv preprint arXiv:1707.07012 (2017).