# Distributed Deep Learning

Amir H. Payberah
payberah@kth.se
2021-12-08
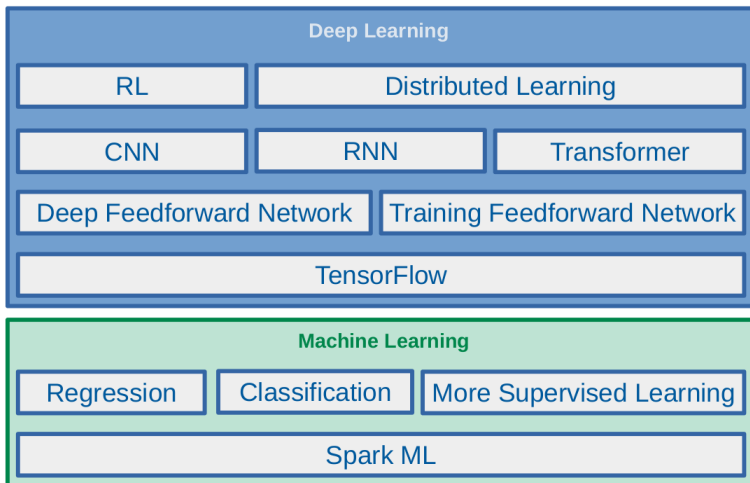
https://id2223kth.github.io

https://tinyurl.com/6s5jy46a

What is the problem?
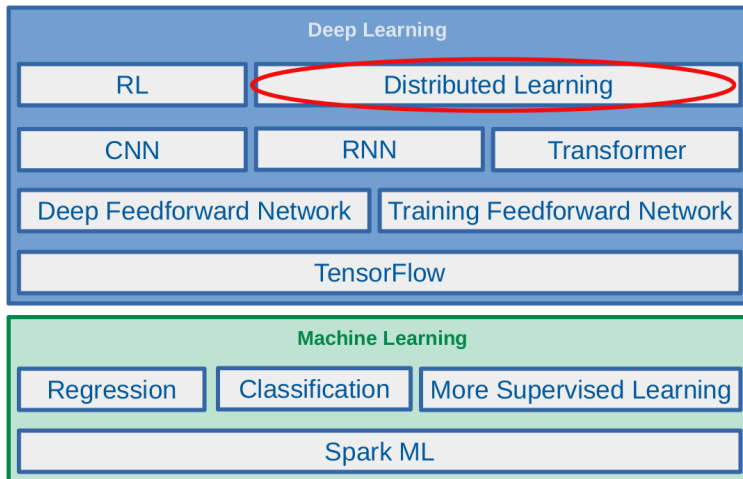
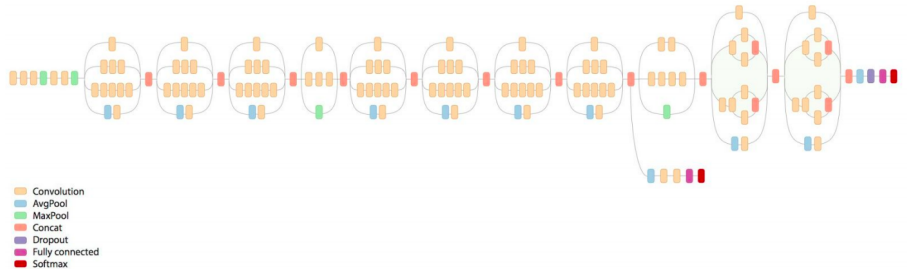- Computationally intensive
- Time consuming



[https://cloud.google.com/tpu/docs/images/inceptionv3onc--oview.png]
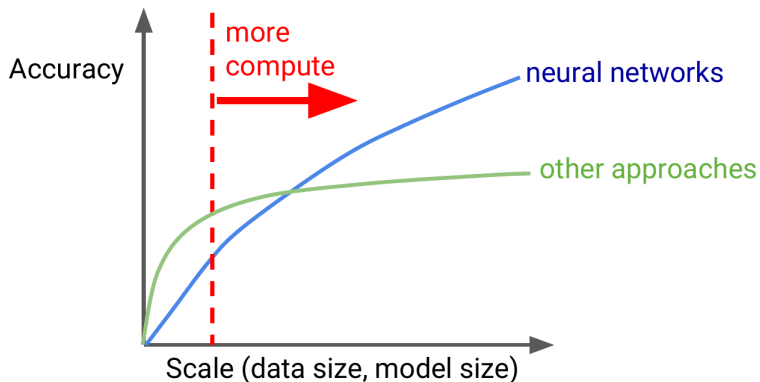
- Massive amount of training dataset
- Large number of parameters

**1980s and 1990s**



[Jeff Dean at AI Frontiers:  Trends and Developments in Deep Learning Research]

**1980s and 1990s**



[Jeff Dean at AI Frontiers:  Trends and Developments in Deep Learning Research]

# Now



[Jeff Dean at AI Frontiers:  Trends and Developments in Deep Learning Research]

# Fundamentals of Machine Learning

# Training Dataset

- E.g., tabular data, image, text, etc.

- E.g., linear models, neural networks, etc.
- $\hat{y} = f_w(\mathbf{x})$

# Loss function

▶ How good $\hat{y}$ is able to predict the expected outcome $y$.

▶ $J(w) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$



▶ E.g., $J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$

# Objective

- <span style="color:green">Minimize</span> the loss function

- $\arg\min_w J(w)$

- $J(w) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$

- $J(w) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$

- Gradient descent, i.e., $w := w - \eta \nabla J(w)$

- Stochastic gradient descent, i.e., $w := w - \eta \tilde{g} J(w)$
  - $\tilde{g}$: gradient at a randomly chosen point.



- Mini-barch gradient descent, i.e., $w := w - \eta \tilde{g}_B J(w)$
  - $\tilde{g}$: gradient with respect to a set of B randomly chosen points.

# Let's Scale the Learning

- Data parallelism

- Model parallelism

# Data Parallelism

# Data Parallelization (1/4)

- Replicate a whole model on every device.
- Train all replicas simultaneously, using a different mini-batch for each.



[Tang et al., Communication-Efficient Distributed Deep Learning:  A Comprehensive Survey, 2020]

# Data Parallelization (2/4)

- $\mathtt{k}$ devices

- $\mathtt{J_j(w)} = \sum_{\mathtt{i}=1}^{\mathtt{b_j}} \mathtt{l(y_i, \hat{y}_i)}, \ \forall \mathtt{j} = 1, 2, \cdots, \mathtt{k}$

- $\tilde{\mathtt{g}}_\mathtt{B}\mathtt{J_j(w)}$: gradient of $\mathtt{J_j(w)}$ with respect to a set of $\mathtt{B}$ randomly chosen points at device $\mathtt{j}$.

- Compute $\tilde{\mathtt{g}}_\mathtt{B}\mathtt{J_j(w)}$ on each device $\mathtt{j}$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Compute the mean of the gradients.
- $\tilde{g}_B J(w) = \frac{1}{k} \sum_{j=1}^{k} \tilde{g}_B J_j(w)$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Update the model.
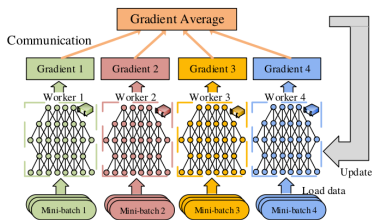- $w := w - \eta \tilde{g}_B J(w)$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

# Data Parallelization Design Issues

- The aggregation algorithm

- Communication synchronization and frequency

- Communication compression

# The Aggregation Algorithm

# The Aggregation Algorithm

- How to aggregate gradients (compute the mean of the gradients)?

- Centralized - parameter server

- Decentralized - all-reduce

- Decentralized - gossip

# Aggregation - Centralized - Parameter Server

- Store the model parameters outside of the workers.

- Workers periodically report their computed parameters or parameter updates to a (set of) parameter server(s) (PSs).



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Mirror all the model parameters across all workers (no PS).
- Workers exchange parameter updates directly via an allreduce operation.



[Tang et al., Communication-Efficient Distributed Deep Learning:  A Comprehensive Survey, 2020]

- No PS, and no global model.
- Every worker communicates updates with their neighbors.
- The consistency of parameters across all workers only at the end of the algorithm.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Reduce: reducing a set of numbers into a smaller set of numbers via a function.
- E.g., sum([1, 2, 3, 4, 5]) = 15
- Reduce takes an array of input elements on each process and returns an array of output elements to the root process.



[https://mpitutorial.com/tutorials/mpi-reduce-and-allreduce]

- AllReduce stores reduced results across all processes rather than the root process.



[https://mpitutorial.com/tutorials/mpi-reduce-and-allreduce]

Initial state

After AllReduce operation

[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- All-to-all allreduce

- Master-worker allreduce

- Tree allreduce

- Round-robin allreduce

- Butterfly allreduce

- Ring allreduce

- **Send** the array of data to each other.
- Apply the reduction operation on each process.
- **Too many** unnecessary messages.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

# AllReduce Implementation - Master-Worker AllReduce

▶ Selecting one process as a master, gather all arrays into the master.
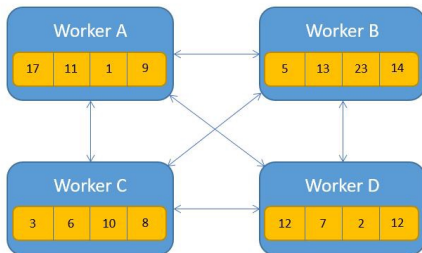
▶ Perform reduction operations locally in the master.

▶ Distribute the result to the other processes.

▶ The master becomes a bottleneck (not scalable).



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

▶ Some try to minimize bandwidth.

▶ Some try to minimize latency.



(a) Tree AllReduce　　　(b) Round-robin AllReduce　　　(c) Butterfly AllReduce

[Zhao H. et al., arXiv:1312.3020, 2013]

- The Ring-Allreduce has two phases:
  1. First, the share-reduce phase
  2. Then, the share-only phase

- In the share-reduce phase, each process `p` sends data to the process `(p+1)%m`
  - `m` is the number of processes, and `%` is the modulo operator.
- The array of data on each process is divided to `m` chunks (`m=4` here).
- Each one of these chunks will be indexed by `i` going forward.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- In the first share-reduce step, process A sends $a_0$ to process B.
- Process B sends $b_1$ to process C, etc.
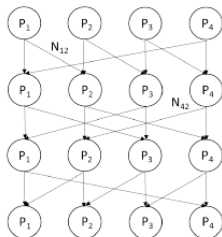


[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- When each process receives the data from the previous process, it applies the reduce operator (e.g., sum)
  - The reduce operator should be associative and commutative.
- It then proceeds to send it to the next process in the ring.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

▶ The share-reduce phase finishes when each process holds the complete reduction of chunk i.

▶ At this point each process holds a part of the end result.



$r_i = a_i + b_i + c_i + d_i$

[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- The share-only step is the same process of sharing the data in a ring-like fashion without applying the reduce operation.

- This consolidates the result of each chunk in every process.



$r_i = a_i + b_i + c_i + d_i$

[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- $N$: number of elements, $m$: number of processes

- Master-Worker AllReduce
  - First each process sends $N$ elements to the master: $N \times (m-1)$ messages.
  - Then the master sends the results back to the process: another $N \times (m-1)$ messages.
  - Total network traffic is $2(N \times (m-1))$, which is proportional to $m$.

- Ring-AllReduce
  - In the share-reduce step each process sends $\frac{N}{m}$ elements, and it does it $m-1$ times: $\frac{N}{m} \times (m-1)$ messages.
  - On the share-only step, each process sends the result for the chunk it calculated: another $\frac{N}{m} \times (m-1)$ messages.
  - Total network traffic is $2(\frac{N}{m} \times (m-1))$.

# Communication Synchronization and Frequency

# Synchronization

- When to synchronize the parameters among the parallel workers?

# Communication Synchronization (1/2)

- Synchronizing the model replicas in data-parallel training requires communication
  - between workers, in allreduce
  - between workers and parameter servers, in the centralized architecture

- The communication synchronization decides how frequently all local models are synchronized with others.

# Communication Synchronization (2/2)

- It will influence:
  - The communication traffic
  - The performance
  - The convergence of model training

- There is a trade-off between the communication traffic and the convergence.

- Two directions for improvement:

  1. To relax the synchronization among all workers.

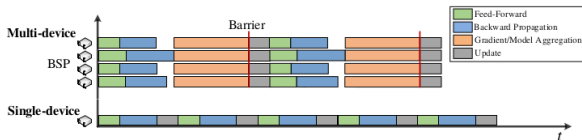  2. The frequency of communication can be reduced by more computation in one iteration.

- Synchronous

- Stale-synchronous

- Asynchronous

- Local SGD

# Communication Synchronization - Synchronous
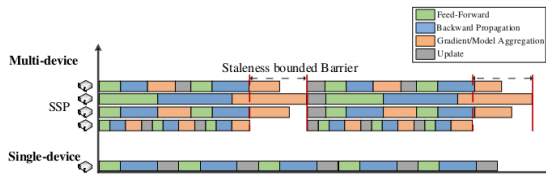
▶ After each iteration, the workers synchronize their parameter updates.

▶ Every worker must wait for all workers to finish the transmission of all parameters in the current iteration, before the next training.

▶ Stragglers can influence the overall system throughput.

▶ High communication cost that limits the system scalability.



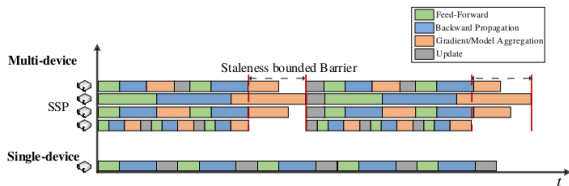[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Alleviate the straggler problem without losing synchronization.

- The faster workers to do more updates than the slower workers to reduce the waiting time of the faster workers.

- Staleness bounded barrier to limit the iteration gap between the fastest worker and the slowest worker.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]
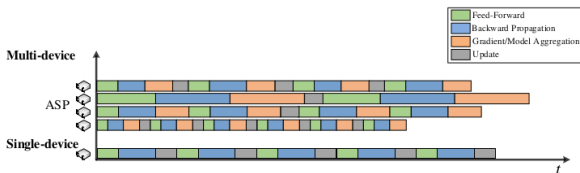
- For a maximum staleness bound $s$, the update formula of worker $i$ at iteration $t+1$:

- $w_{i,t+1} := w_0 - \eta(\sum_{k=1}^{t}\sum_{j=1}^{n} G_{j,k} + \sum_{k=t-s}^{t} G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$

- The update has three parts:
  1. Guaranteed pre-window updates from clock $1$ to $t$ over all workers.
  2. Guaranteed read-my-writes in-window updates made by the querying worker $i$.
  3. Best-effort in-window updates. $S_{i,t+1}$ is some subset of the updates from other workers during period $[t-s]$.



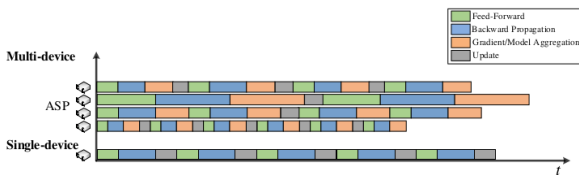[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- It completely eliminates the synchronization.

- Each work transmits its gradients to the PS after it calculates the gradients.

- The PS updates the global model without waiting for the other workers.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

▶ $w_{t+1} := w_t - \eta \sum_{i=1}^{n} G_{i,t-\tau_{k,i}}$

▶ $\tau_{k,i}$ is the time delay between the moment when worker $i$ calculates the gradient at the current iteration.
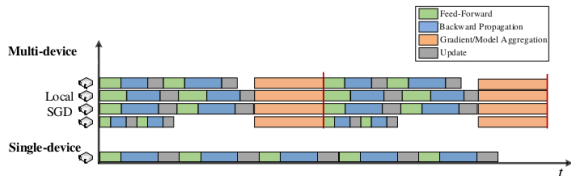


[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

▶ All workers run several iterations, and then averages all local models into the newest global model.

▶ If $\mathcal{I}_T$ represents the synchronization timestamps, then:

$$w_{i,t+1} = \begin{cases} w_{i,t} - \eta G_{i,t} & \text{if} \quad t+1 \notin \mathcal{I}_T \\ w_{i,t} - \eta \frac{1}{n} \sum_{i=1}^{n} G_{i,t} & \text{if} \quad t+1 \in \mathcal{I}_T \end{cases}$$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]
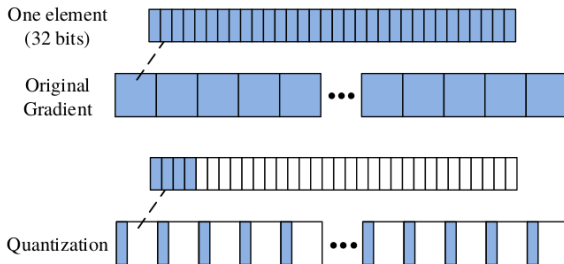
# Communication Compression

# Communication Compression

- Reduce the communication traffic with little impact on the model convergence.

- Compress the exchanged gradients or models before transmitting across the network.

- Quantization

- Sparsification

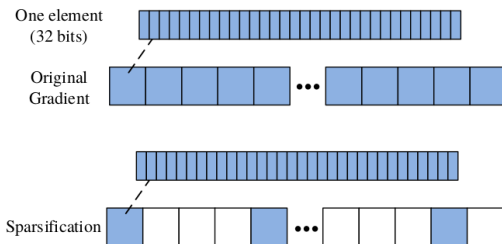- Useing lower bits to represent the data.

- The gradients are of low precision.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Reducing the number of elements that are transmitted at each iteration.

- Only significant gradients are required to update the model parameter to guarantee the convergence of the training.

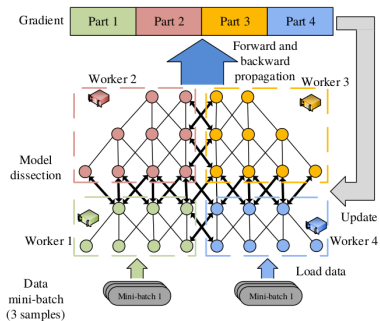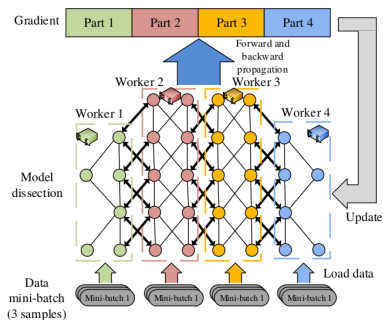- E.g., the zero-valued elements are no need to transmit.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

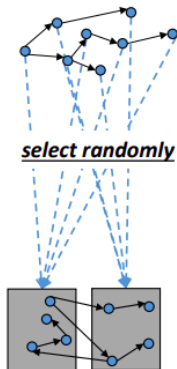# Model Parallelism

# Model Parallelization

▶ The model is split across multiple devices.

▶ Depends on the architecture of the NN.



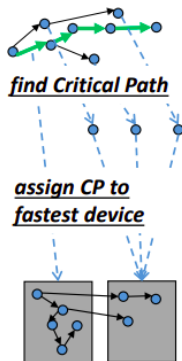[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Randomly assign vertices to devices proportionally to the capacity of the devices by using a hash function.


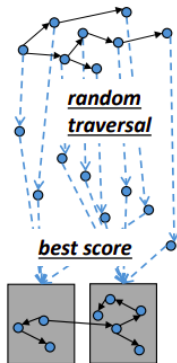
[Mayer, R. et al., The TensorFlow Partitioning and Scheduling Problem, 2017]

- Assigning the complete critical path to the fastest device.
- Critical path: the path with the longest computation time from source to sink vertex.



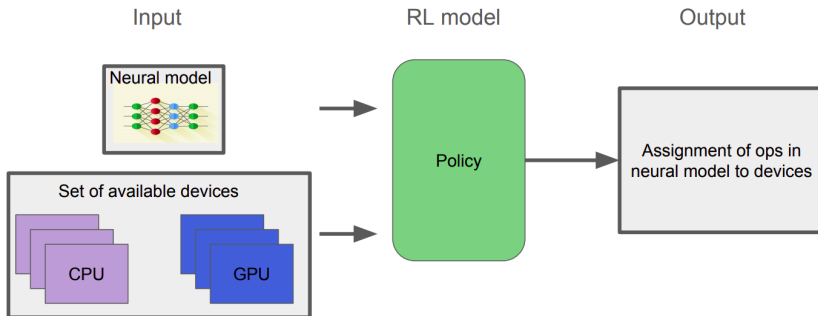[Mayer, R. et al., The TensorFlow Partitioning and Scheduling Problem, 2017]

- Different objectives, e.g., memory, importance, traffic, and execution time



[Mayer, R. et al., The TensorFlow Partitioning and Scheduling Problem, 2017]

[Mirhoseini et al., Device Placement Optimization with Reinforcement Learning, 2017]

[Mirhoseini et al., Device Placement Optimization with Reinforcement Learning, 2017]

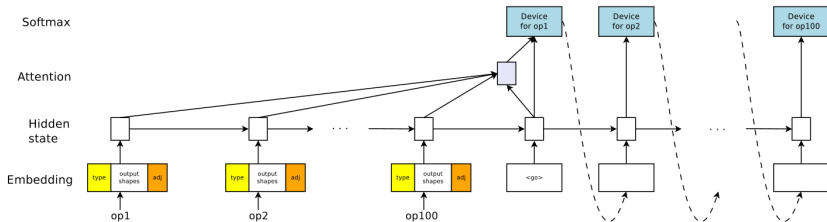- $J(w) = \mathbb{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G}, w)}[R(\mathcal{P})|\mathcal{G}]$
- Objective: $\arg\min_w J(w)$

- $\mathcal{G}$: input neural graph
- $R$: runtime
- $J(w)$: expected runtime
- $w$: trainable parameters of policy
- $\pi(\mathcal{P}|\mathcal{G}, w)$: policy
- $\mathcal{P}$: output placements $\in \{1, 2, ..., num\_ops\}^{num\_devices}$
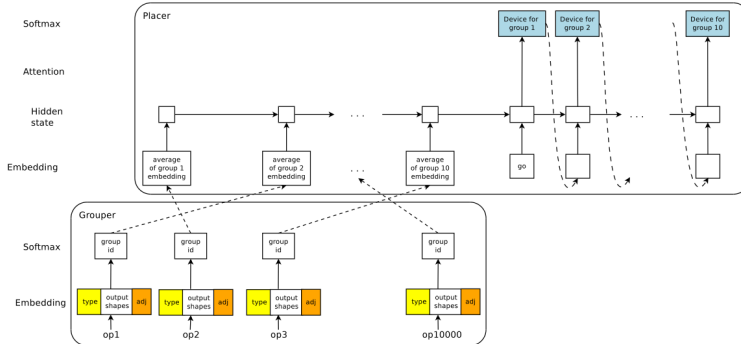
- RL reward function based on execution runtime.
- The RL policy is defined as a seq-to-seq model.
- RNN Encoder receives graph embedding for each operation.
- RNN Decoder predicts a device placement for each operation.



[Mirhoseini et al., Device Placement Optimization with Reinforcement Learning, 2017]

- Grouping operations.
- Prediction is for group placement, not for a single operation.



[Mirhoseini et al., A Hierarchical Model for Device Placement, 2018]

# Summary

# Summary

- ▶ Scalability matters

- ▶ Parallelization

- ▶ Data Parallelization
  - Parameter server vs. AllReduce
  - Synchronized vs. asynchronized

- ▶ Model Parallelization
  - Random, critical path, multi-objective, RL

Thanks!