# Scalable Distributed Machine Learning with Huge Data for IoT and Scientific Discovery: Opportunities and Challenges

Tony Luo        Sajal Das

Department of Computer Science

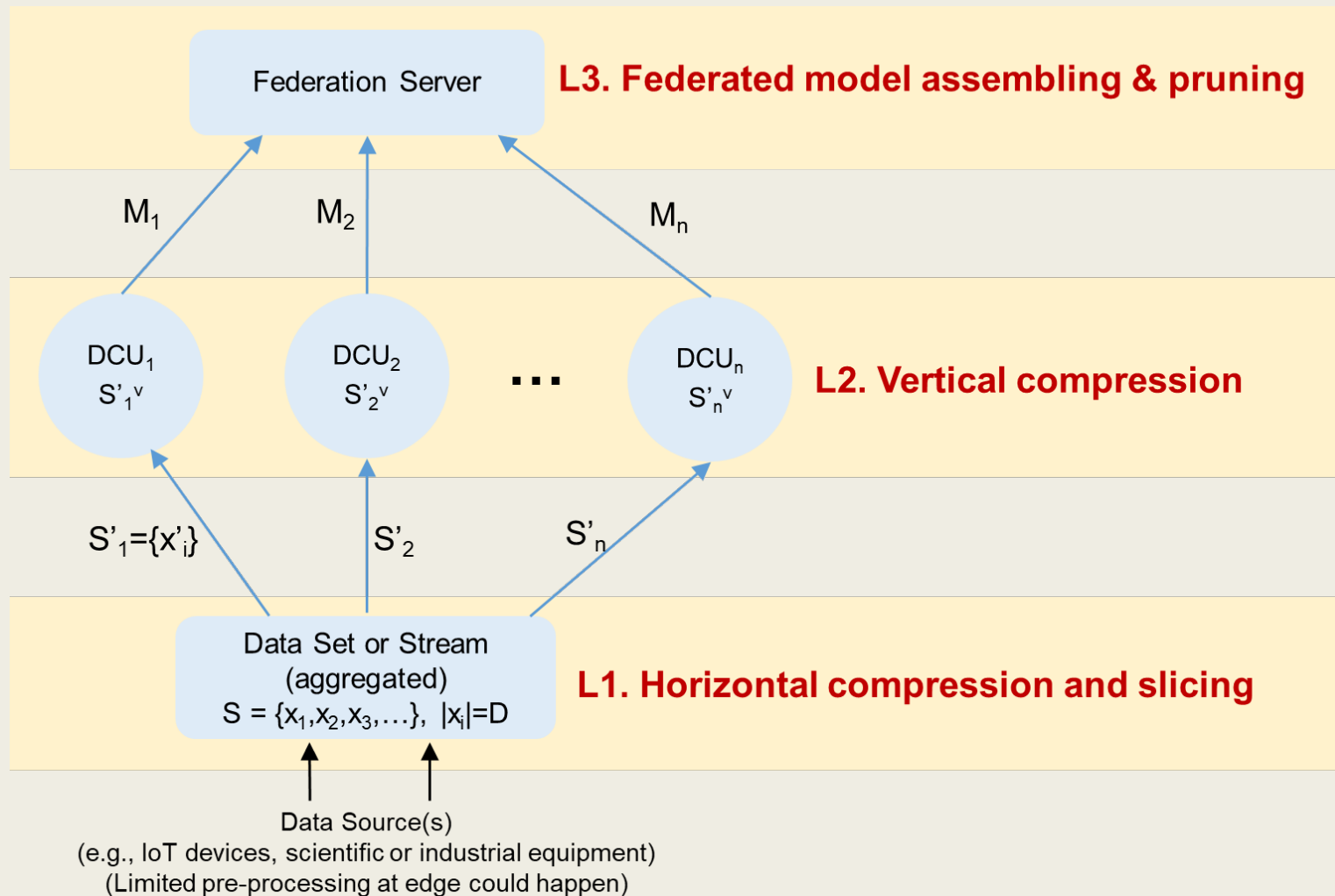Missouri University of Science and Technology

14 April 2020

# Motivation and Objective

■ "Huge Data" problem: Data generation rates in petabyte, exabyte, and even zettabyte per hour are becoming increasingly common

– Scientific fields such as astronomy, physics, and biological sciences

– Internet of Things (IoT) with billions of sensors and devices interconnected



■ Our goal: solve the "Huge Data" problem in the context of machine learning such that models can be efficiently trained on a huge amount of data and can make fast predictions on new data samples.
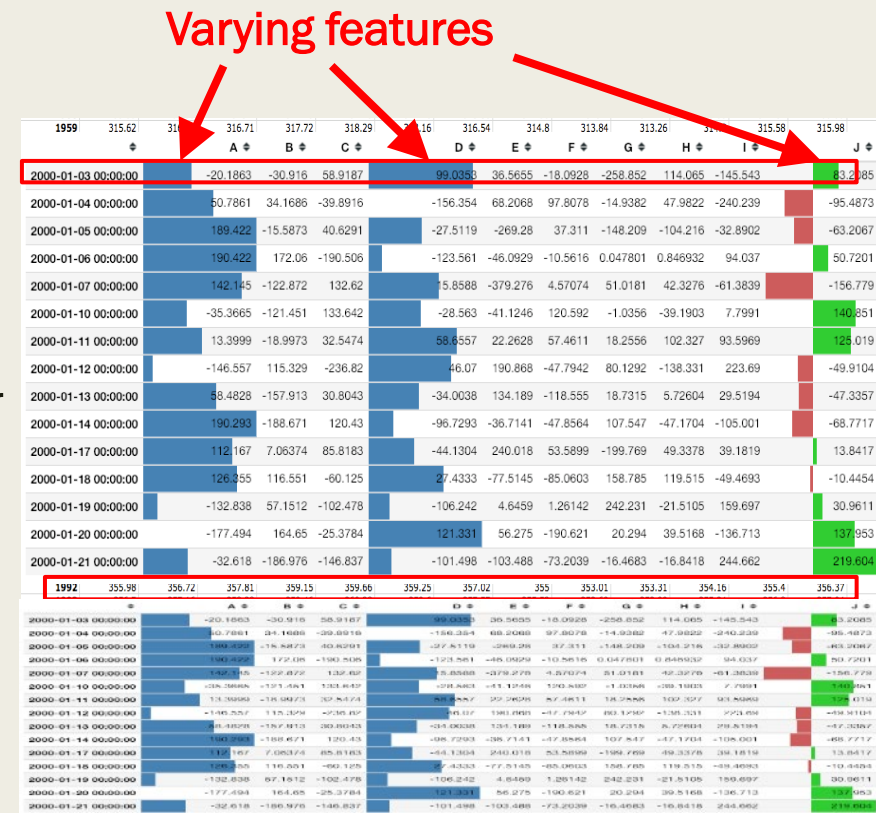
# Scalable distributed ML framework



- Original data set or stream
  - $S = \{x1, x2, x3, ...\}$

- Huge in two dimensions:
  - **Vertically huge:** a huge number of data samples (i.e., rows) or the influx streaming rate is extremely high
  - **Horizontally huge:** a huge number of features (i.e., columns), D
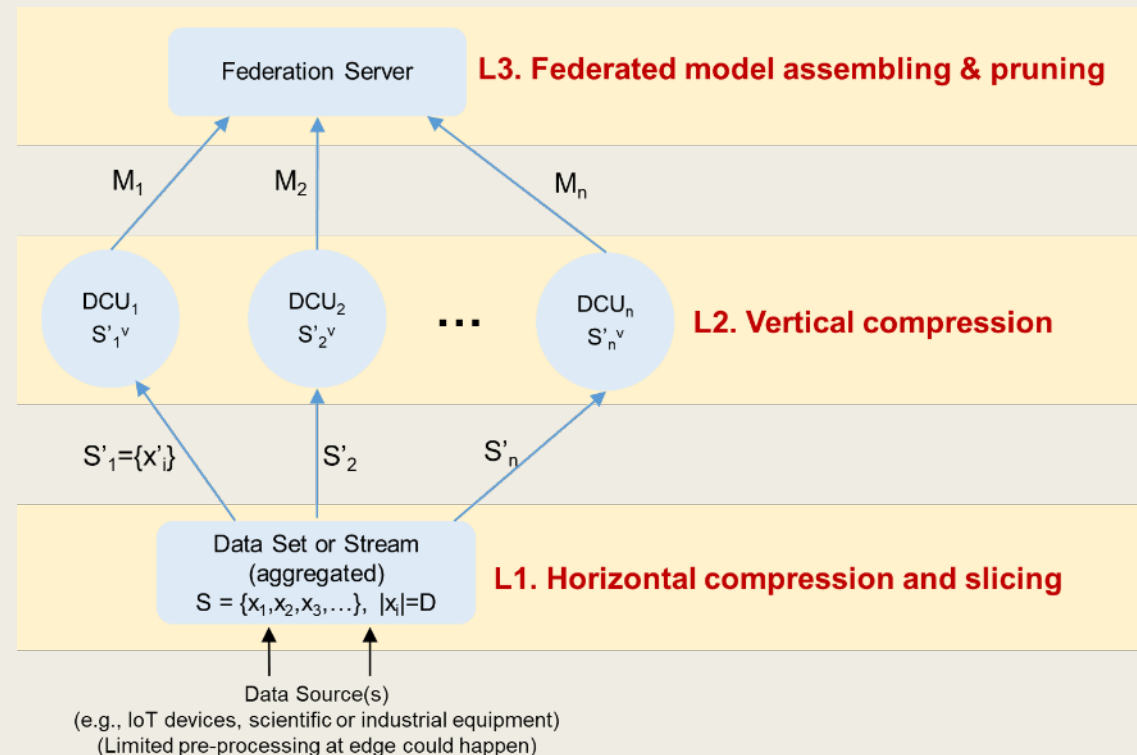
# L1. Horizontal compression and slicing

- **Horizontal compression**: reduce the number of rows

- Key observation: consecutive data samples often have little variation except for a few features, whose variation can be characterized by a well-defined function or probability distribution

- Compression method: only keep a "seed sample" + descriptive functions; all the subsequent samples are drastically condensed into a single number, until another "seed sample" appears, which indicates a change that the descriptive functions cannot describe.

  – Original S → Seed sample set S'={x'1, x'2, x'3,…} + metadata set Q = {{F1, k1}, {F2, k2}, {F3, k3},…}, $k_j$ is the number of data samples described by tuple $(x'_j, F_j)$

  – Compression ratio: $r = 1 - m/\sum k_j$ where m=|S'| is the number of signal samples

  – $r \to 1$ in many practical cases!



Varying features

- **<u>Horizontal slicing</u>**: slice the compressed data S' horizontally, into $S'_1, S'_2,..., S'n$ while keeping intact the width of each data sample D. The subsets $S'_1, S'_2,..., S'n$ are then sent to *n* distributed computing units (DCUs)

- When S is a data stream rather than a static dataset, the slicing procedure is substituted by a dispatching mechanism
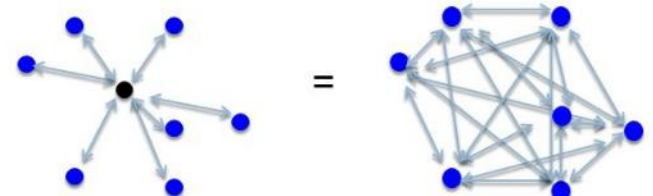
Does not wait for each subset $S'i$ to be complete but immediately sends it to a DCU once the size of $S'i$ reaches a certain value (can be as small as one), and this repeats.

Choice of DCU: round-robin or adaptively based on bandwith and workload.



Federation Server — **L3. Federated model assembling & pruning**

$M_1$ $M_2$ $M_n$

$DCU_1$ $S'^v_1$  $DCU_2$ $S'^v_2$ ... $DCU_n$ $S'^v_n$ — **L2. Vertical compression**

$S'_1=\{x'_i\}$ $S'_2$ $S'_n$

Data Set or Stream (aggregated)
$S = \{x_1, x_2, x_3,...\}$, $|x_i|=D$ — **L1. Horizontal compression and slicing**

Data Source(s)
(e.g., IoT devices, scientific or industrial equipment)
(Limited pre-processing at edge could happen)

# L2. Vertical compression

■ At each DCU: reduce the dimension of each sample from D to d (d<<D)

■ Classic dimensionality reduction technique: principal component analysis (PCA)

 – Deterministic: always result in the same subset of reduced features

 – Not desired by some ML tasks such as ensemble learning, which needs model diversity to boost performance

■ We propose random Johnson-Lindenstrauss projection as an alternative

 – Introduces randomness while preserving pairwise distances

$$\sum_{i=1}^{n} \left\| \mathbf{a}_i - \mu_{C(i)} \right\|_2^2 = \sum_{i=1}^{k} \frac{1}{|C_i|} \sum_{(w,v)\in C_i} \left\| \mathbf{a}_w - \mathbf{a}_v \right\|_2^2$$
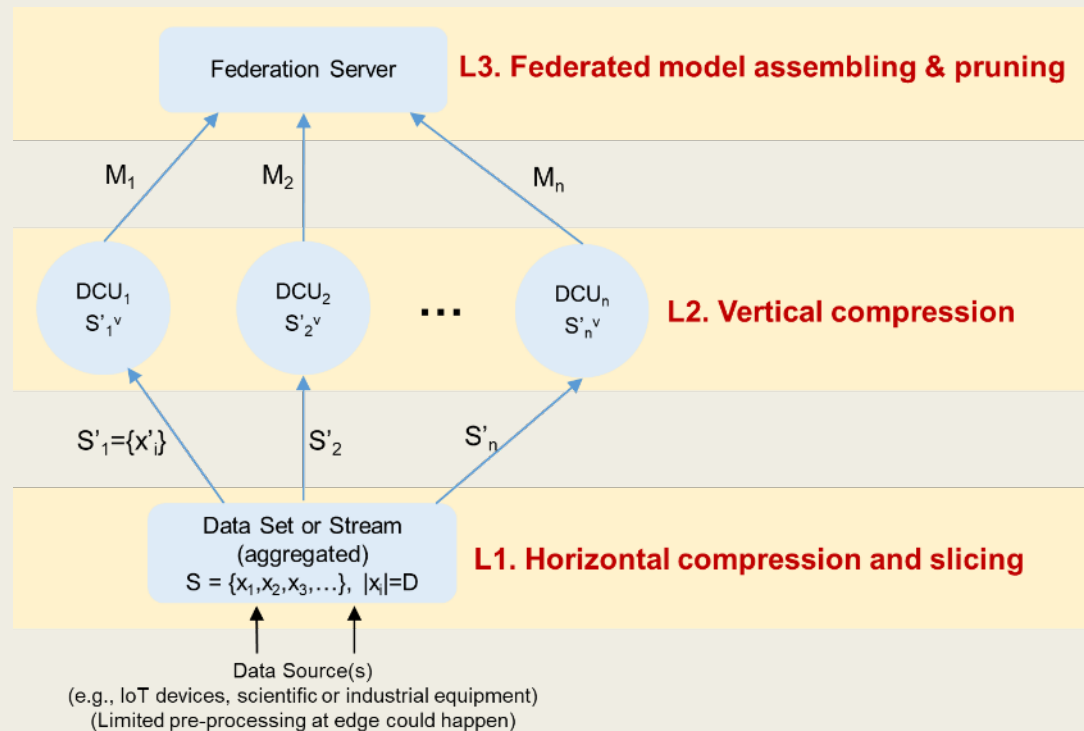


■ After vertically compressing S'$_i$ into S'$_i^v$, a sub-model M$_i$ is trained over S'$_i^v$ on the same DCU

 – Training process needs to be tailored to accommodate the metadata introduced by L1, for which we anticipate a "virtual expansion" procedure

# L3.  Federated model assembling & pruning

- **Assembling**: aggregate the n sub-models $M_i$ into one global model M

- **Challenge**: sub-models are heterogeneous
  - Horizontal slicing (L1) may result in non-i.i.d. datasets
  - Vertical compression (L2) may result in different features
  - Conventional distributed machine learning not applicable
  - Federated averaging can handle non-iid data but not feature heterogeneity
  - A novel method is needed: Open challenge

- **Pruning**: global model M may still contain a large number of parameters due to the huge data setting, and thus may need a pruning process before deployment in order to achieve fast runtime predictions
  - Starting point: exist pruning and quantization techniques for deep neural networks (DNN)

# Conclusion

■ Proposed a framework and preliminary ideas to advance the state-of-the-art in scalable distributed machine learning for huge data analytics

■ Our approach could alleviate the hurdles of storing, transferring, and processing huge data such that scientific research and IoT systems at huge-data scale could substantially benefit from AI



Tony Luo:
tluo@mst.edu
https://tluocs.github.io