

TRUONG THAO NGUYEN

産業技術総合研究所実社会ビッグデータ活用オープンイノベーションラボラトリ
Postdoctoral Researcher

分散型ディープニューラルネットワークの大規模設計の調査・研究

§ 1. 研究成果の概要

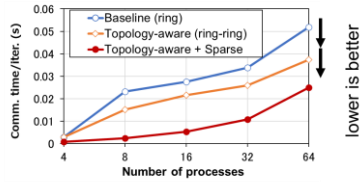
The target of this project is to enable training Deep Learning model on large-scale HPC system with a short time including study (1) new parallelism strategies (not only data parallelism) (2) the method to reduce the communication time and (3) the effect of different system architecture on training time. Based on this research, we aim to develop an estimating model as the basis for a utility we name it ParaDL, that aids end-users, framework developers and system builders in (i) identifying the optimal large-scale parallelism strategies, i.e. domain decomposition, when deploying a training job for a given DL model and data set on an HPC system; (ii) guiding the implementation and possible optimizations of different parallel strategies in existing DL frameworks; and (iii) advising system architects on the best co-design choices for their system depending on the workloads they plan to run

Task 1: We formally define the possible basis parallelism strategies and build up an analysis and estimation model for these types of strategies. We collaborated with Barcelona Supercomputing Center (BSC) to compare the performance projected by my model with the those implemented in a real system, e.g., ABCI (implementation by BSC). We submit this work to HPDC2020 (A* conference). We plan to find a method for searching the best (fine-grained) hybrid parallelism, i.e., applying different parallelism approach for different layers of a DNN model in the future.

Topology-aware Allreduce
 ✓ Reduce comm. time up to 45%
 ✓ Reduce power consumption. up to 23%



Approximate Sparse data
 ✓ 100x-1000x compressed
 ✓ Reduce communication time ~40% more



Simulated result with ABCI-system, 32MB-message, 0.78% sparcification

Task 2: Basically, state of the art techniques for optimizing communication time are: (i) Architecture-aware algorithm, (ii) data compression and (iii) comm. /comp. overlapping. We proposed Allreduce optimized for GPU-cluster (CCPE journal, I/F 1.167) and its combination with direction (ii) (IPCCC conference). This helps to reduce the communication time up to 60%. In the future, we plan to find the best combination of those three techniques.

Task 3: We did the survey on the trends of system architecture supported for Deep Learning. In the next fiscal year, we plan to estimate our model on different system (not only ABCI). This work may require collaboration with another research center.