

ONNX Model Zoo

Distributed Training

Presenter(s): Rodolfo Esteves, Rajeev Nalawadi
Intel Corporation

Model Zoo Training samples

With inclusion of ONNX operators for training

- ONNX Model zoo should cover few model samples in zoo for training scenarios
- Target at least 1 model to begin. Could be any category Vision, NLP etc..
 - Preferably Transformer/NLP training sample as first priority
 - Request is to cover a model in each category (3 model samples would provide enough starting momentum to showcase similarities/differences)
- ONNX will be able to show distributed training infrastructure by leveraging (Horovod, Deepspeed/ZERO, ...) across clusters and the industry SOTA techniques of data/model/pipeline parallelisms (as applicable)
- Quantization aware training (QAT) can be better demonstrated with ONNX model samples (with original model converted from other FW's in FP32 and then quantization flows applied using ONNX)
 - Original FP32 model facilitates model accuracy comparison
- Mixed precisions usage for training the ONNX models can also be highlighted in future

Thank You !!

