

Model	Speedup
GNMT	1.7X speedup
NCF	2.6X speedup
ResNet-50-v1.5	3.3X speedup
SSD-RN50-FPN-640	2.5X speedup

7.2.3. Tensor Core Optimized Model Scripts For TensorFlow

The [tensor core examples provided in GitHub](#) focus on achieving the best performance and convergence from NVIDIA Volta tensor cores by using the latest [deep learning example](#) networks and [model scripts](#) for training.

Each example model trains with mixed precision tensor cores on Volta, therefore you can get results much faster than training without tensor cores. This model is tested against each NGC monthly container release to ensure consistent accuracy and performance over time. The TensorFlow container includes the following [TensorFlow tensor core examples](#):

- ▶ An implementation of the [SSD320 v1.2](#) model. The SSD320 v1.2 model is based on the [SSD: Single Shot MultiBox Detector](#) paper, which describes SSD as “a method for detecting objects in images using a single deep neural network”. Our implementation is based on the existing [model from the TensorFlow models repository](#).
- ▶ An implementation of the [Neural Collaborative Filtering \(NCF\)](#) model. The NCF model is a neural network that provides collaborative filtering based on implicit feedback, specifically, it provides product recommendations based on user and item interactions. The training data for this model should contain a sequence of user ID, item ID pairs indicating that the specified user has interacted with, for example, was given a rating to or clicked on, the specified item.
- ▶ An implementation of the [Bert](#) model. BERT, or Bidirectional Encoder Representations from Transformers, is a new method of pre-training language representations which obtains state-of-the-art results on a wide array of Natural Language Processing (NLP) tasks. This model is based on [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#) paper. NVIDIA's BERT is an optimized version of [Google's official implementation](#), leveraging mixed precision arithmetic and tensor cores on V100 GPUS for faster training times while maintaining target accuracy.
- ▶ An implementation of the [U-Net Industrial Defect Segmentation](#) model. This U-Net model is adapted from the original version of the [U-Net model](#) which is a convolutional auto-encoder for 2D image segmentation. U-Net was first introduced by Olaf Ronneberger, Philip Fischer, and Thomas Brox in the paper: [U-Net: Convolutional Networks for Biomedical Image Segmentation](#). This work proposes a modified version of U-Net, called TinyUNet which performs efficiently and with very high accuracy on the industrial anomaly dataset [DAGM2007](#).