

📖 Note Information

Author : AbbasXu

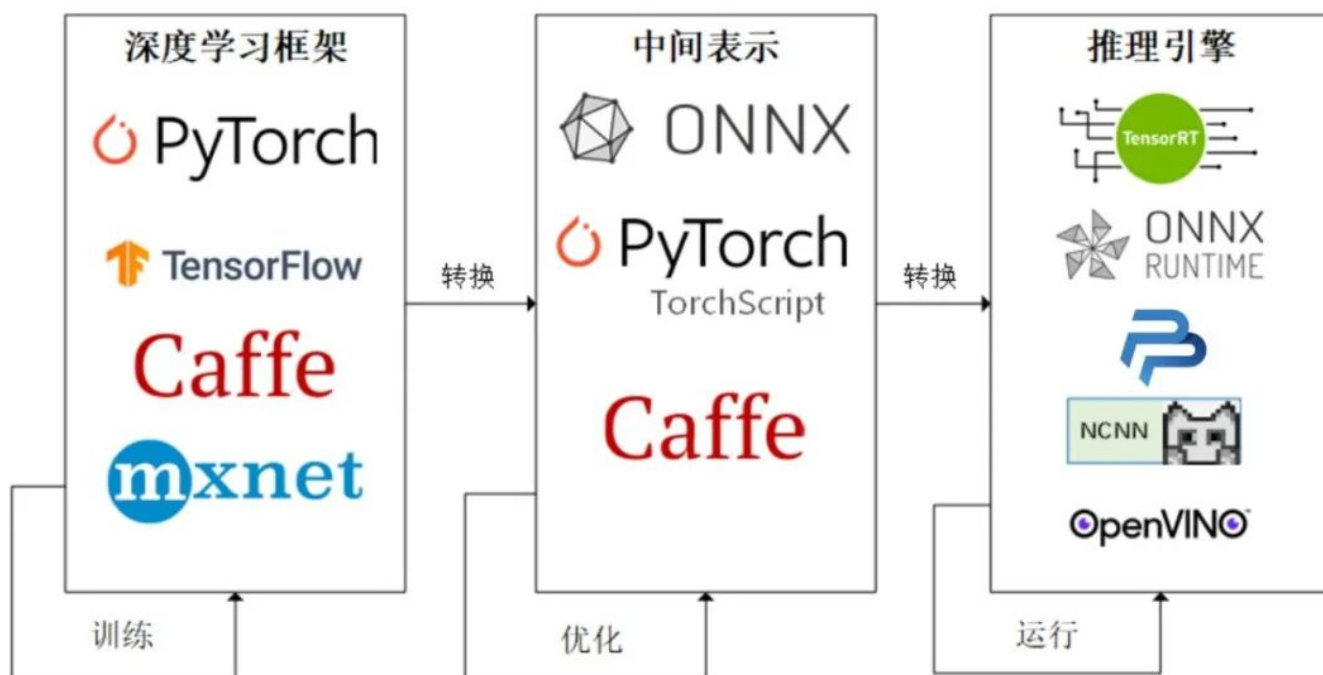
Date : 2022-08-28

Title : pytorch基础

Keywords : #pytorch #模型部署

使用ONNX进行部署并推理

模型部署pipeline



ONNX简介

- ONNX官网: <https://onnx.ai/>
- ONNX GitHub: <https://github.com/onnx/onnx>

ONNX(Open Neural Network Exchange) 是 Facebook (现Meta) 和微软在2017年共同发布的, 用于标准描述计算图的一种格式。

ONNX 已经对接了下图的多种深度学习框架和多种推理引擎。

Frameworks & Converters

Use the frameworks you already know and love.



Transformers



Keras

LibSVM

MATLAB®

[M]^s MindSpore



NCNN



PaddlePaddle



SIEMENS



SciKit Learn



ONNX Runtime简介

- ONNX Runtime官网: <https://www.onnxruntime.ai/>
- ONNX Runtime GitHub: <https://github.com/microsoft/onnxruntime>
ONNX Runtime 是由微软维护的一个跨平台机器学习推理加速器，它直接对接ONNX，可以直接读取.onnx文件并实现推理，不需要再把 .onnx 格式的文件转换成其他格式的文件。PyTorch借助ONNX Runtime也完成了部署的最后一公里，构建了 PyTorch --> ONNX --> ONNX Runtime 部署流水线，我们只需要将模型转换为 .onnx 文件，并在 ONNX Runtime 上运行模型即可。

模型导出为ONNX

使用 `torch.onnx.export()` 把模型转换成 ONNX 格式的函数。

通过 `onnx.checker.check_model()` 进行检验。

使用 [Netron](#) ONNX可视化

