

ONNX

(Open Neural Network Exchange)

Current status & beyond

Motivation: Lots of DL frameworks!



CNTK



And more!

- Create open standard format for deep learning models
- Provide deep learning model interoperability among frameworks
- Protobuf: github.com/onnx/onnx/blob/master/onnx/onnx.proto

Industry wide participation

At inception:



Microsoft

Dec 2017

Industry partners:



Current Spec as of May 2019:

1.5

Current status: Framework & converter support

8 Frameworks



```
#Export Alexnet from PyTorch
import torch
torch.onnx.export(model, dummy_input, "alexnet.onnx", ...)

#Import Alexnet to Caffe2
import onnx
import caffe2.python.onnx.backend as backend
model = onnx.load("alexnet.onnx")
rep = backend.prepare(model, device="CPU")
```

6 Converters



```
#To convert TF model to ONNX model
python -m tf2onnx.convert\
    --input tests/models/fc-layers/frozen.pb\
    --output tests/models/fc-layers/model.onnx\
```

Current status: Runtimes, Compilers, & Visualizers

11 Runtimes



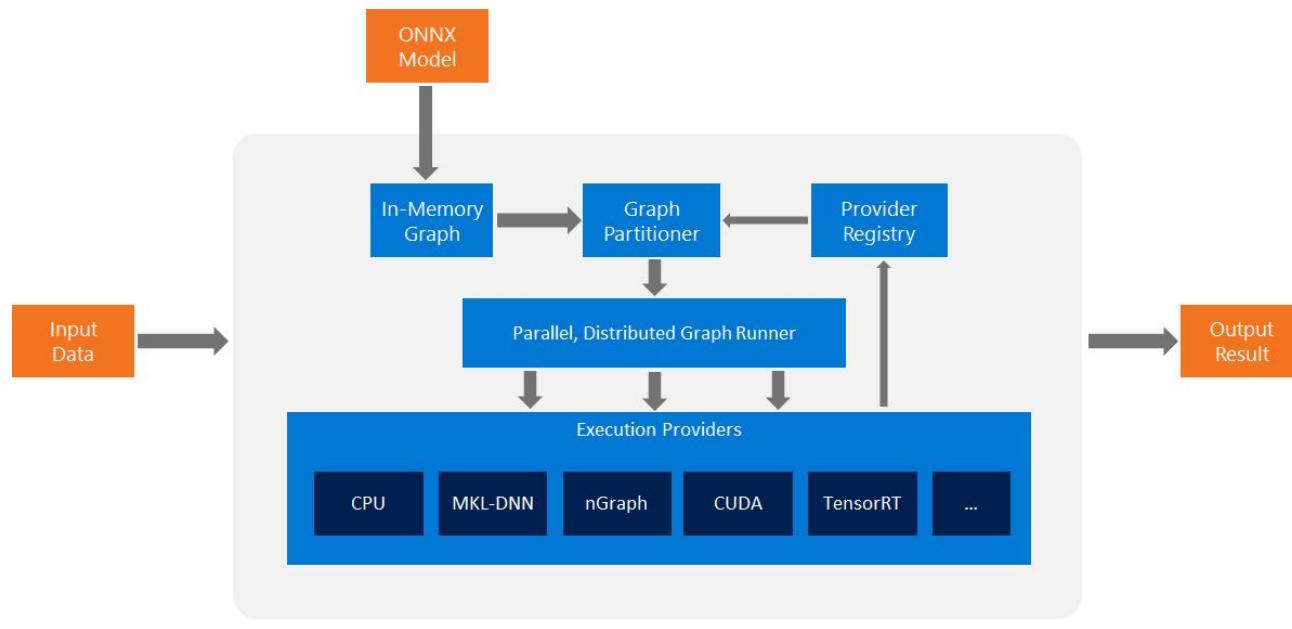
4 Compilers



2 Visualizers



Current status: ONNX Runtime



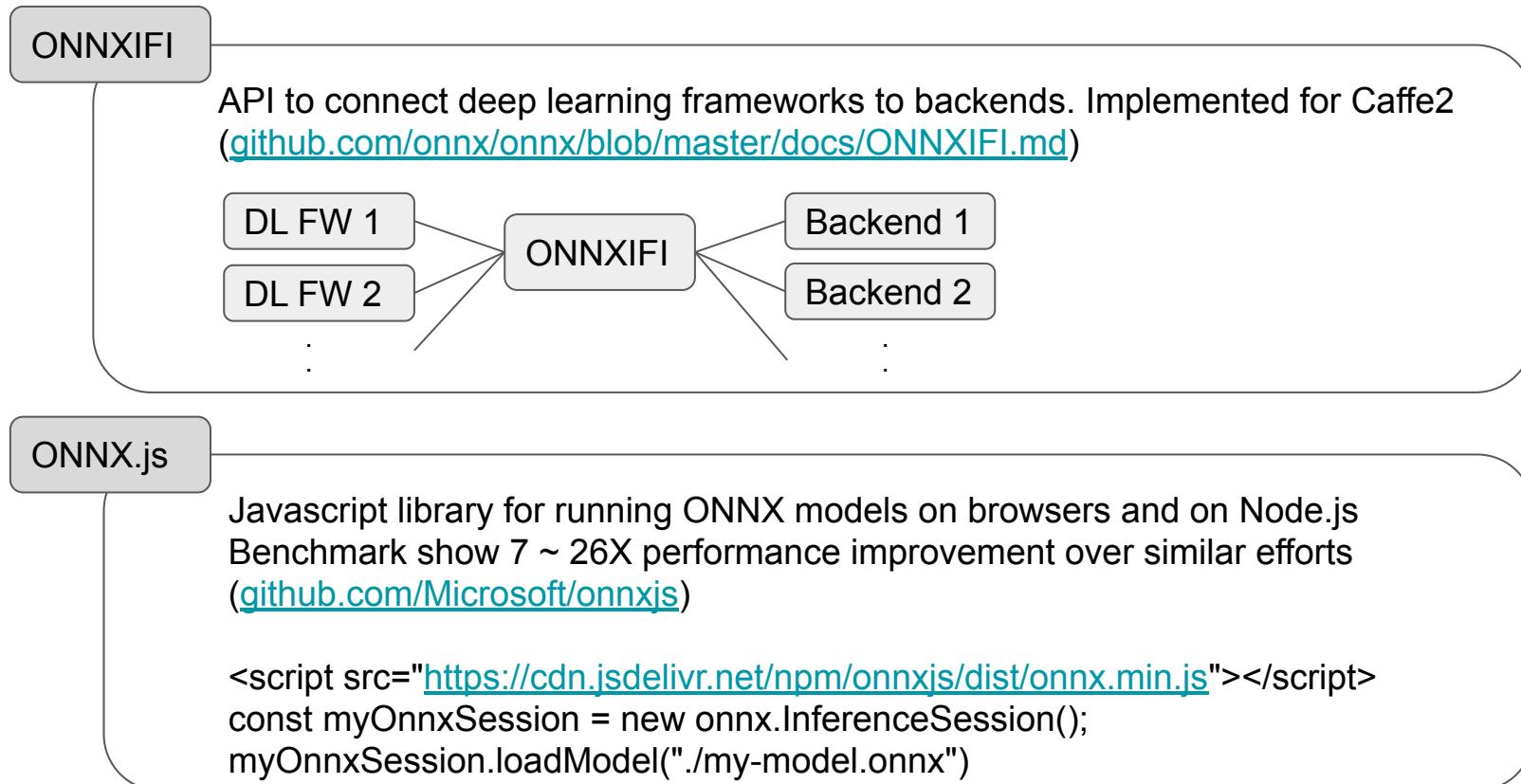
Will be open sourced at //BUILD conference (May 6)

On average 2X performance boost*

Supports quantization

More EPs in progress

Current status: ONNXIFI & ONNX.js



Current status: Model zoo (github.com/ONNX/models)

Image Classification

MobileNet	Bvlc_reference_CaffeNet
ResNet	Bvlc_reference_RCNN_ILSVRC13
SqueezeNet	DenseNet121
VGG	Inception_v1
Bvlc_AlexNet	Inception_v2
Bvlc_GoogleNet	ShuffleNet
	ZFNet512

Object Detection & Segmentation

Tiny_YOLOv2
SSD **New!**
YOLO v3 **New!**
DUC

Body, Face & Gesture Analysis

ArcFace
Emotion FerPlus

CNN models

Building blocks for RNN available

Current status: ONNX 1.5 spec

- Opset 10 adds operators to support object detection models such as Yolo v3, Faster RCNN, and SSD. Models will be added to the ONNX Model Zoo
- Quantization support (with first set of operators)
- Promote ONNX Function to support composing operators (support of more operators from other frameworks while limiting new operators)
- All experimental ops are removed and deprecated

Areas of development (gitter.im/onnx/Lobby)

- **Quantization:** Few more quantized ops to be added in the future (Intel)
- **Training:** Use cases defined. Technical discussion on enabling training IR, loss function & optimizer. Enabling ONNX Runtime for training (IBM)
- **Model zoo:** Implementing CI to automate compliance check (Huawei)
- **Mobile & Edge:** defining profiles for edge and associated use cases (Qualcomm)
- **Dataflow:** Define preprocessing op definitions. Reference implementation to be provided (NVidia)

New governance structure for wider participation

Transparent decision process & better technical decisions

github.com/onnx/onnx/tree/master/community

Community

Members

Contributors

Approvers

Member
companies

Structure

Steering Committee

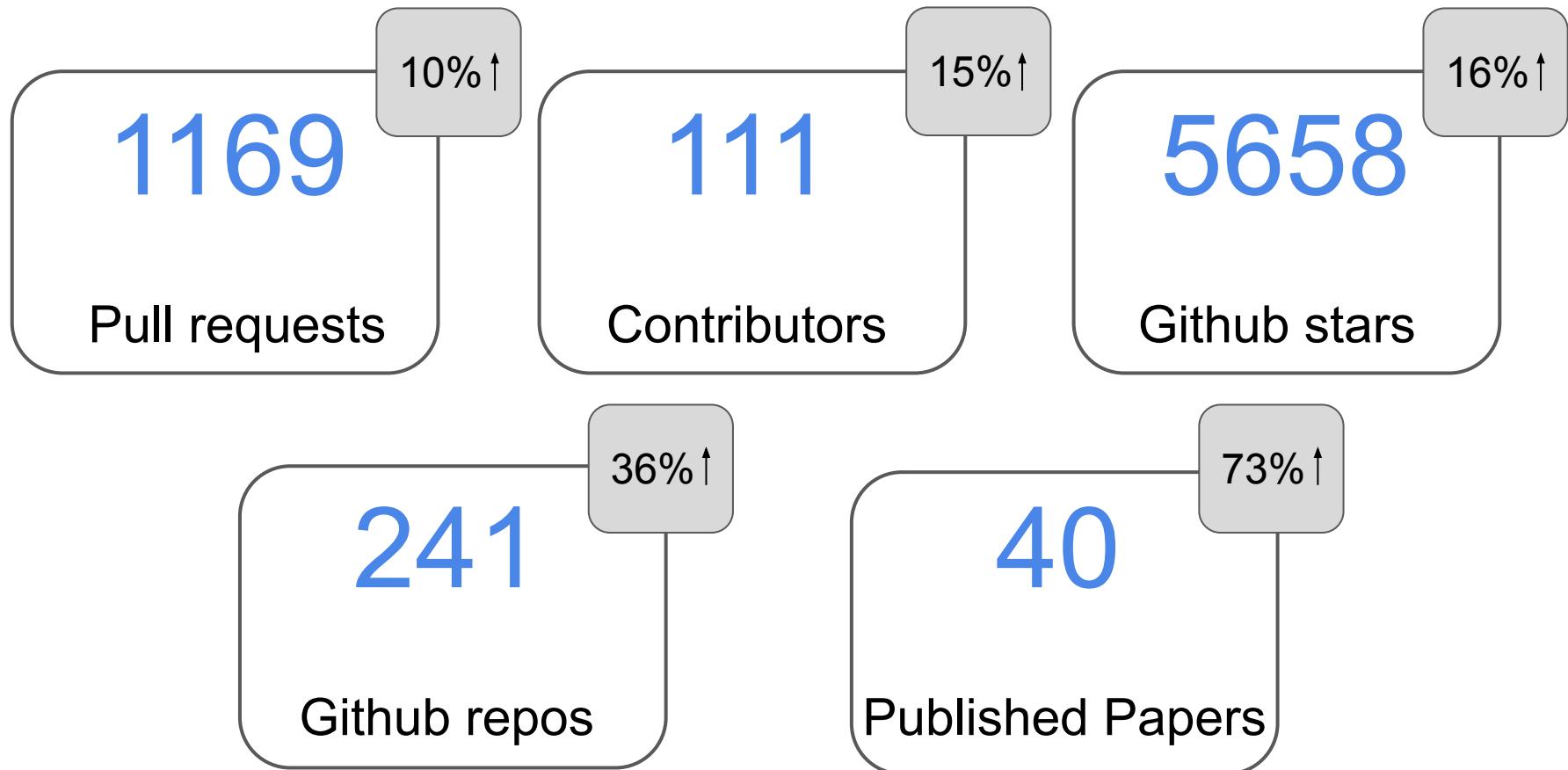
Special Interest Groups

Working groups

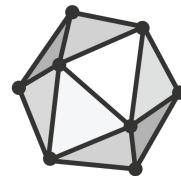
Persistent

Temporary

Engagement & usage (Dec 2018 => Mar 2019)



Thank you!



onnx.ai