# Update

- 2 releases since last meet-up
  - 1.13 (Przemyslaw Wysocki, Intel)
  - 1.14 (Yuan Yao, Nvidia)
- IR 9
  - 8-bit Float
  - More support to local functions
- Reference implementation for ONNX operators

# Update (Continued)

- Serialization with Textproto format
- Function Ops can have multiple bodies, each uses a different opset version
- Inline Local functions
- Documentation page: Comparing between versions of an operator

# Highlight Op changes on the New Documentation page

# Validate PyTorch Converter with ONNX Reference Implementation

```
sess = onnxreference.ReferenceEvaluator(onnx_model)

onnx_out = sess.run(None, input_dict)
```
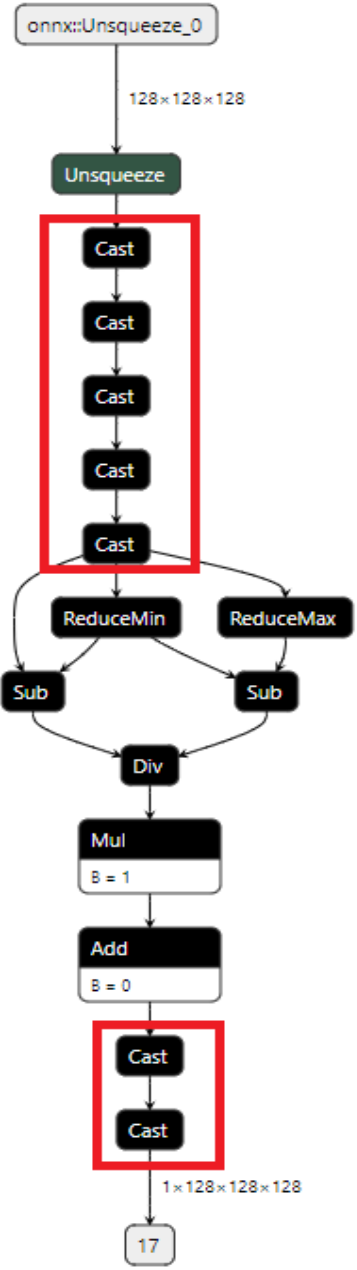
Code    Blame    1113 lines (941 loc) · 46 KB    Raw

```python
584    def convert_to_onnx(
661                torch.onnx.export(
662                    mode_to_export,
663                    tuple(inputs),
664                    f=filename,
665                    input_names=input_names,
666                    output_names=output_names,
667                    dynamic_axes=dynamic_axes,
668                    opset_version=opset_version,
669                    **torch_versioned_kwargs,
670                )
671                onnx_model = onnx.load(filename)
672
673            if verify:
674                if device is None:
675                    device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
676
677                inputs = [i.to(device) if isinstance(i, torch.Tensor) else i for i in inputs]
678                model = model.to(device)
679
680                with torch.no_grad():
681                    set_determinism(seed=0)
682                    torch_out = ensure_tuple(model(*inputs), True)
683
684                set_determinism(seed=0)
685                model_input_names = [i.name for i in onnx_model.graph.input]
686                input_dict = dict(zip(model_input_names, [i.cpu().numpy() for i in inputs]))
687                if use_ort:
688                    ort_sess = onnxruntime.InferenceSession(
689                        onnx_model.SerializeToString(), providers=ort_provider if ort_provider else ["CPUExecutionProvider"]
690                    )
691                    onnx_out = ort_sess.run(None, input_dict)
692                else:
693                    sess = onnxreference.ReferenceEvaluator(onnx_model)
694                    onnx_out = sess.run(None, input_dict)
695                set_determinism(seed=None)
696                # compare onnx/ort and PyTorch results
697                for r1, r2 in zip(torch_out, onnx_out):
698                    if isinstance(r1, torch.Tensor):
699                        assert_fn = torch.testing.assert_close if pytorch_after(1, 11) else torch.testing.assert_allclose
700                        assert_fn(r1.cpu(), convert_to_tensor(r2, dtype=r1.dtype), rtol=rtol, atol=atol)  # type: ignore
701
702            return onnx_model
703
```

Use Text Editor To Edit An ONNX Model

# Roadmap

- Maintain a stable CI Infrastructure

- Translate more operators to function op

- Work with operator, pre/post processing, and converter SIGs to support real world models, in production

# ONNX with Project MONAI

- [https://github.com/Project-MONAI](https://github.com/Project-MONAI)
- Home to most PyTorch models in medical imaging domain
- 20+ SOTA models published in Model-zoo
- ONNX utility to convert all models to ONNX format
- Utilize ORT and ONNX reference implementation for model validation
- However, running core ONNX models is only part of medical imaging workflow
- Need to put end-to-end MONAI workflow in ONNX so it can be applied anywhere on any hardware

Medical Imaging Workflow

# ONNX Pre and Post Processing

- Pre- and Post-Processing operations in MONAI are performed using the "Compose" class, which encapsulates transform sequences.

- These operations are wrapped into a Torch module and exported.

- The 'export_compose' function handles exporting the Compose object to ONNX

```
def export_compose(pnp_compose: Compose, opset_version:
int, inputs: Sequence[np.ndarray], outputs:
Sequence[np.ndarray], image_meta_dict: Dict[str, Any],
task_name: str) -> ModelProto:
```

- If new Ops need to be added, we will make it function op. (For example: AffineGrid.)

# ONNX Inferer

- Work with SlidingWindowInferer first – it is used by majority of models in model-zoo

- With ONNX-Script, the op can be easily implemented as a function op.

- Combining with ONNX PNP, we complete medical imaging workflow in ONNX.

- Any models developed with MONAI framework can be converted to ONNX automatically

```python
@script()
def sliding_window_inference(inputs: FLOAT["N", "C", "D", "H", "W"], roi_size: INT64[3]) -> FLOAT["N", "Seg_C", "D", "H",
"W"]:
    """
    The sliding window method is used for model inference. It involves taking a 3D sliding window on the input tensor
    and making predictions using a provided predictor. The outputs from the predictor are then aggregated to form the
output of the operator.
    """
    inputs_shape = op.Shape(inputs)
    inputs_spatial_shape = op.Shape(inputs, from=2)
    N, _, D, H, W = op.Split(inputs_shape, num_outputs=5)
    roi_D, roi_H, roi_W = op.Split(roi_size, num_outputs=3)

    scan_interval = roi_size
    slices = dense_patch_slices_script(inputs_spatial_shape, roi_size, scan_interval)
    S_, _, _ = op.Split(op.Shape(slices), num_outputs=3)
    S = op.Squeeze(S_, op.Constant(value_ints=[0]))

    seg_C = op.Constant(value_ints=[2]) # TODO: get from predictor model
    output_shape = op.Concat(N, seg_C, inputs_spatial_shape, axis=0)

    aggrregated_pred = op.CastLike(op.ConstantOfShape(output_shape), inputs)
    aggrregated_count = op.CastLike(op.ConstantOfShape(inputs_shape), roi_size)
    for slice_g in range(S):
        win_data, start, stop = prepare_for_predictor_batch_size_is_1_script(inputs, slice_g, slices)
        pred = op.OpaqueOp(win_data, model_path="C:/Temp/sliding_window_predictor_sw_batch_size_is_1.onnx")
        aggrregated_pred, aggrregated_count = aggrregate_predictor_output(pred, start, stop, aggrregated_pred,
aggrregated_count)

    return aggrregated_pred / op.CastLike(aggrregated_count, aggrregated_pred)
```
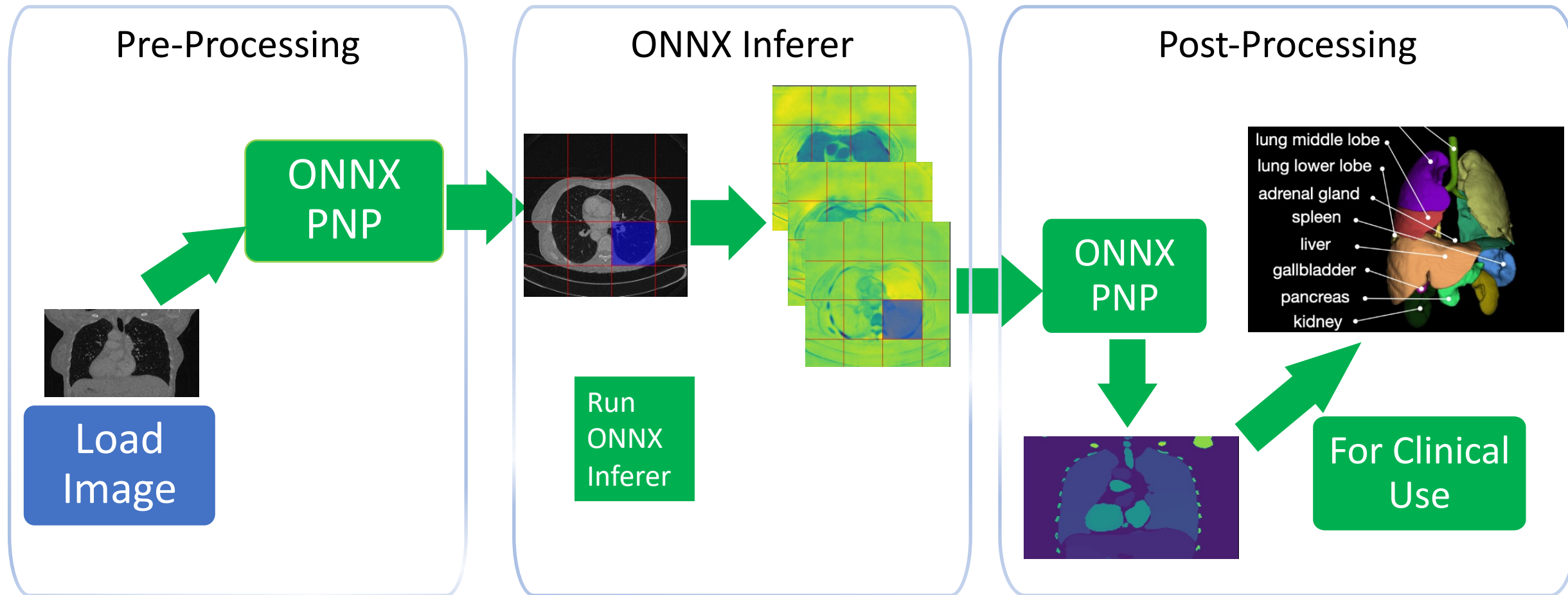
# Complete Medical Imaging Workflow with ONNX

Please Get Involved!

Github: PRs, Issues, and Discussions

Slack channel: https://slack.lfai.foundation and join onnx-archinfra

Monthly SIG meetings (see slack channel for announcements)

Thanks you!

# Useful links

- [Comparing ONNX operator Split - 13 vs 18](#)
- [Load](#) and [Save](#) ONNX models in textproto format
- [ONNX exporter in MONAI. It uses ONNX reference implementation to validate the converted models](#)
- [Affine](#), [GridSample](#) PRs for image pre- and post-processing
- [SlidingWindowInferer](#)