



<https://tinyurl.com/tract-rs>

slides

tract
an efficient rust neural network
inference engine
&
Shipping neural networks with
torch-to-nnef

Overview

- Mathieu and Julien, ML software engineers for Sonos Voice Control
- Introducing two open-source libraries
 - tract, a Rust generic neural network inference library
 - torch-to-nnef, a companion project to export PyTorch model in tract's preferred format





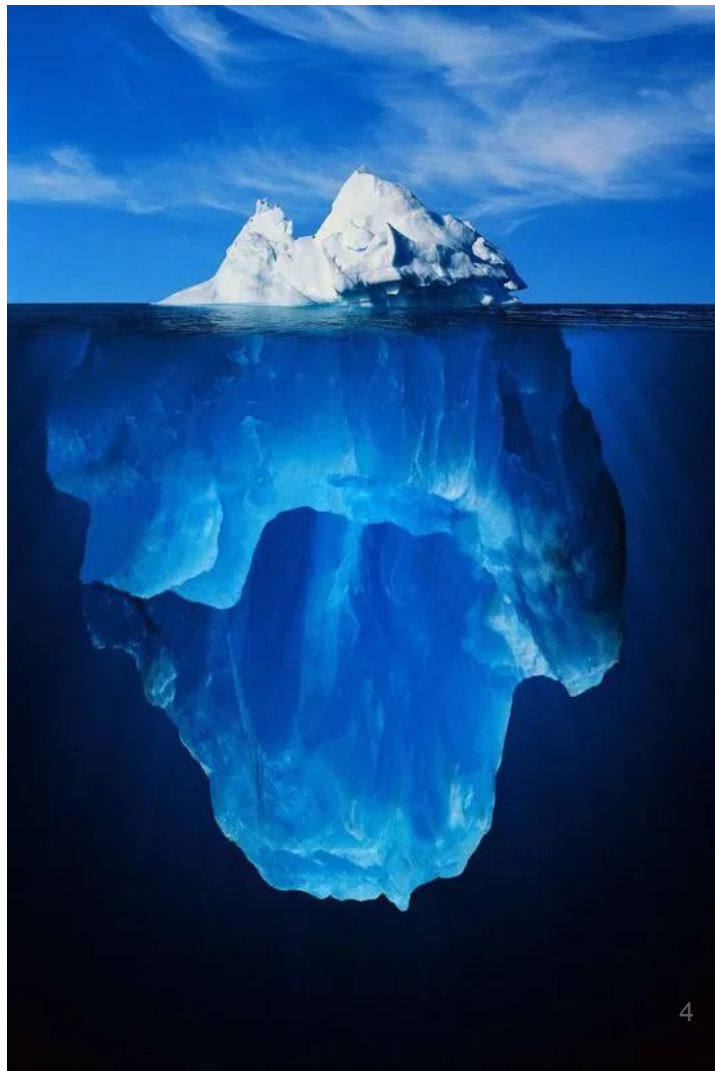
tract

an efficient rust neural network inference engine

15min

Scratching the surface

- Neural Network inference library
- “tract all the way down”: full stack, from model parsing to optimized kernels
- Have been around for 9 years or so
- In Rust, integration friendly, Apache2/MIT



```

fn main() -> Result<()> {
    let nnef = tract_rs::nnef()?;
    let model = nnef.load("mobilenet_v2_1.0.onnx.nnef.tgz)?.into_runnable()?;

    // open image, resize it and make a Tensor out of it
    let image = image::open("grace_hopper.jpg").unwrap().to_rgb8();
    let resized =
        image::imageops::resize(&image, 224, 224, ::image::imageops::FilterType::Triangle);
    let input = tract_ndarray::Array4::from_shape_fn((1, 3, 224, 224), |(_, c, y, x)| {
        let mean = [0.485, 0.456, 0.406][c];
        let std = [0.229, 0.224, 0.225][c];
        (resized[(x as _, y as _)][c] as f32 / 255.0 - mean) / std
    }));

    // run the model on the input
    let result: Vec<Value> = model.run([input])?;

    // find and display the max value with its index
    let best =
        &result[0].as_slice::<f32>()?.iter().zip(2..).max_by(|a, b| a.0.partial_cmp(b.0).unwrap());
    println!("result: {best:?}");
    Ok(())
}

```

tract mobilenet_v2_1.0.nnetf.tgz

```
0 Source data
  — 1,3,224,224,F32
11 3 Conv conv
  * Data format: NCHW
  * Kernel shape:[3, 3] (strides:Some([2, 2]), padding:Explicit([1, 1], [1, 1]), dilations:Some([1, 1]))
  * Kernel OIHW (groups:1)
  — 1,32,112,112,F32
5 5 Max relu_relu_y_4
7 7 EinSum conv_1
  * NIHW,OI->OHW (F32)
  — 32,112,112,F32
9 9 Sub batch_normalization_1_batch_normalization_output_3
11 11 Mul batch_normalization_1_batch_normalization_output_8
13 13 Mul batch_normalization_1_batch_normalization_output_16
```


tract mobilenet_v2_1.0.nnetf.tgz -O dump --cost --profile -R

0.000 ms/i	0.0%				0 Source data
					— 1,3,224,224,F32
0.033 ms/i	0.1%				1 Pad conv.pad
					— 1,3,226,226,F32
0.001 ms/i	0.0%				2 LazyIm2col conv.lazyIm2col
					— 1,1,Opaque DynPackedOpaqueFact { k: Val(27), mn: Val(12544), packers: [Pa
0.273 ms/i	0.7%	FMA(F32)	10838016	39.717 GF/s	5 OptMatMul conv.matmatmul
					— 1,1,32,12544,F32
0.027 ms/i	0.1%				6 IntoShape conv.reshape_group
					— 1,32,112,112,F32
0.140 ms/i	0.3%				8 OptMaxByScalar relu_relu_y_4
0.065 ms/i	0.2%				9 OptMatMulPack conv_1.pack_a
					— 1,112,Opaque DynPackedOpaqueFact { k: Val(32), mn: Val(112), packers: [Pa
0.310 ms/i	0.8%	FMA(F32)	12845056	41.401 GF/s	16 OptMatMul relu_1_relu_y_0
					— 32,112,112,F32
3.692 ms/i	9.1%	FMA(F32)	3612672	978.558 MF/s	19 DepthWiseConv conv_2
0.085 ms/i	0.2%				21 OptSubByScalar batch_normalization_2_batch_normalization_output_3
0.084 ms/i	0.2%	FMA(F32)	401408	4.751 GF/s	23 OptMulByScalar batch_normalization_2_batch_normalization_output_8
0.084 ms/i	0.2%	FMA(F32)	401408	4.795 GF/s	25 OptMulByScalar batch_normalization_2_batch_normalization_output_16
0.084 ms/i	0.2%				27 OptAddByScalar batch_normalization_2_batch_normalization_output_21
0.138 ms/i	0.3%				28 OptMaxByScalar relu_2_relu_y_0
0.213 ms/i	0.5%				29 OptMatMulPack conv_3.pack_a
					— 112,Opaque DynPackedOpaqueFact { k: Val(32), mn: Val(112), packers: [Pack
0.176 ms/i	0.4%	FMA(F32)	6422528	36.474 GF/s	35 OptMatMul batch_normalization_3_batch_normalization_output_21
					— 16,112,112,F32
0.045 ms/i	0.1%				36 OptMatMulPack conv_4.pack_a
					— 112,Opaque DynPackedOpaqueFact { k: Val(16), mn: Val(112), packers: [Pack
0.561 ms/i	1.4%	FMA(F32)	19267584	34.350 GF/s	43 OptMatMul relu_3_relu_y_0
					— 96,112,112,F32
2.748 ms/i	6.8%	FMA(F32)	2709504	986.030 MF/s	46 DepthWiseConv conv_5
					— 96,56,56,F32

Easy to integrate

- No cumbersome third party dependency
- Easy to cross compile (Rust...)
- WASM support (run your model in a browser)
- Optimized for ARM32, ARM64, X64, WASM SIMD, metal, cuda.
- C API
- Python bindings: `pip install tract`



GPU support

```
let nnef = tract_rs::nnef()?;  
let model = nnef.load("mobilenet_v2_1.0.onnx.nnef.tgz")?.into_runnable()?;
```

tract mobilenet_v2_1.0.nnef.tgz -O bench

Bench ran 127 times, 39.493 ms/i.

```
let nnef = tract_rs::nnef()?;  
let model = nnef.load("mobilenet_v2_1.0.onnx.nnef.tgz")?;  
let model = runtime_for_name("cuda")?.prepare(model)?;
```

tract mobilenet_v2_1.0.nnef.tgz --cuda bench

Bench ran 2259 times, 2.204 ms/i.

```
tract --llm --op1 Llama-3.1-8B-Instruct-f16f16.nnef.tgz
```

```
1 Source input_ids
  — 1,S,I64
2 Gather model_model_embedTokens_inputsEmbeds_0
  — 1,S,2048,F16
3 Cast model_model__0_inputLayernorm_hiddenStatesTo0
  — 1,S,2048,F32
4 RmsNorm model_model__0_inputLayernorm_varianceMean0_mean_reduce_output.rms_norm
5 Cast model_model__0_inputLayernorm_to0
  — 1,S,2048,F16
7 Mul model_model__0_inputLayernorm_hiddenStatesMul0_1
  — 1,S,2048,F16
```

```
29 ApplyRope model_model__0_selfAttn_keyStatesAdd0_1.apply_rope
30 DynamicKeyValueCache in_cache_key_0
   — 1,8,S+P,64,F16
```

ONNX models support

- ONNX is first class citizen.
- tract strong tensor typing : for each tensor we want
 - a fixed rank
 - all dimension known (symbolic expression allowed)
- ONNX protobuf format does not contain detailed tensor shapes
- Two-stage loading: load, add some shape hints, tract will infer the rest
- Still supporting symbolic shape

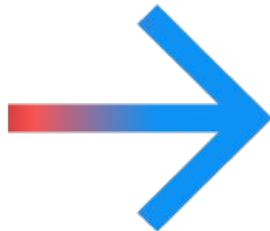
```

#[test]
fn test_inference_model() -> anyhow::Result<()> {
    ensure_models()?;
    let mut model = onnx()?.load("mobilenetv2-7.onnx")?;
    model.set_input_fact(0, "1,3,224,224,F32")?;
    let model = model.into_tract()?.into_runnable()?;
    let hopper = grace_hopper();
    let result = model.run([hopper])?;
    let view = result[0].view::<f32>()?;
    let best = view
        .as_slice()
        .unwrap()
        .iter()
        .enumerate()
        .max_by(|a, b| a.1.partial_cmp(b.1).unwrap())
        .unwrap();
    assert_eq!(best.0, 652);
    Ok(())
}

```

LLMs support and torch-to-nnef

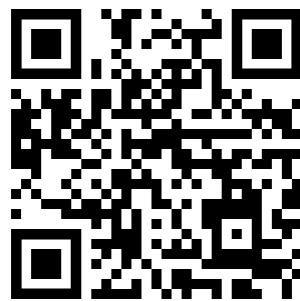
- The LLM tsunami happened with specialized engines
- tract supports LLMs, but as a generic engine
 - Models structures are not hard-coded in tract
 - Instead, tract has an SDPA operator (and a couple other dedicated ones)
 - tract needs to be given the graph description in the input format: NNEF (or ONNX)
- ONNX is a bit awkward with LLM (getting better lately)
- This lead us to open-source torch-to-nnef last summer.



Shipping neural networks with torch-to-nnef

A highly tract compatible exporter

5min



[tinyurl.com/torch-to-nnef
docs](https://tinyurl.com/torch-to-nnef/docs)

What composes a neural network model asset ?

A **neural network asset** is the packaged output of training that is designed for efficient inference. It typically contains two essential components:




A list of **tensors** with names and values (the parameters learned during the training process).



A computation graph stitching the tensors together: inputs, outputs, data types, tensor shapes, and the sequence of operators.

What is NNEF ?

NNEF stands for **N**eural **N**etwork **E**xchange **F**ormat.

It addresses the same core problem as  ONNX but focusing on inference.

Specified by the **K H R  N O S**[®] G R O U P (consortium of ~170 companies)

Why NNEF is interesting ?



Readable Graph Structure

The main **.nnef** file represents the model graph in a **simple, declarative, text-based format**



Composition

Neural-network are made of repetition of blocks, the text format promotes **reusability**, avoids repetition, and enables a **clean functional structure**.



Extensible Tensors & Quantization Logic

- Quantization metadata can live in a **separate .quant** file
- Flexible Tensor Storage**
Each tensor is stored as a binary **.dat** blob.

NNEF syntax example

```
fragment gelu_fast_approx( x: tensor<scalar> ) -> ( y: tensor<scalar> )
{
    y = 0.5 * x * (1.0 + tanh(x * 0.7978845608 * (1.0 + 0.044715 * x * x)));
}

graph network(input_0) -> (output_0)
{
    input_0 = tract_core_external(shape = [1, 2], datum_type = 'f32');
    lin_weight = variable<scalar>(label = 'lin.weight', shape = [2, 2]);
    lin_bias = variable<scalar>(label = 'lin.bias', shape = [2]);
    lin_bias_aligned_rank_expanded = unsqueeze(lin_bias, axes = [0]);
    lin_x_linear0 = linear(input_0, lin_weight, lin_bias_aligned_rank_expanded);
    gelu0 = gelu_fast_approx(lin_x_linear0);
    lin_bias_aligned_rank_expanded = unsqueeze(lin_bias, axes = [0]);
    lin_2nd_call_x = linear(gelu0, lin_weight, lin_bias_aligned_rank_expanded);
    output_0 = gelu_fast_approx(lin_2nd_call_x);
}
```

What is torch-to-nnef ?

Torch to NNEF is a Python package that:

- was born in early 2022 out of the frustration of ONNX quantization export state and our need back then for an internal project
- Creates a strong bridge between PyTorch and NNEF
- Allows a complete and flexible control to export any model.

pip install **torch-to-nnef**

Export Python example

```
from pathlib import Path
from torch_to_nnef import export_model_to_nnef, TractNNEF

export_model_to_nnef(
    # any nn.Module
    model=my_model,
    # tuple of model arguments
    args=inputs_data_sample,
    # filepath to dump NNEF archive
    file_path_export=Path("vit_b_16.nnef.tgz"),
    # inference engine to target
    inference_target=TractNNEF(
        # tract version (to ensure compatible operators)
        version="0.21.13",
        # and correctness of output compared to PyTorch
        # for the provided model and input is performed
        check_io=True,
    ),
    input_names=["input"],
    output_names=["output"],
    # create a debug bundle in case model export work
    # but NNEF fail in tract
    # (either due to load error or precision mismatch)
    debug_bundle_path=Path("./debug.tgz"),
)
```

Why torch-to-nnef ?

- Export advanced quantization/dequantization functions, data-types optimized on PyTorch side seamlessly.
- Optionally control serialization of specific ***nn.Module***
- Provide additional expressions on symbolic dimensions, helping engine to reason about them (ie. Batch dimension being always >0).
- Export signal based networks easily (MFCC, MelFilterBanks, ...)



Thanks for listening

Interested to use or contribute ?



github.com/sonos/tract



github.com/sonos/torch-to-nnef

Documentations and past articles:

- Technical [documentation tract](#)
- Technical [documentation torch-to-nnef](#)
- Sonos tech blog related articles: [\[1\]](#), [\[2\]](#), [\[3\]](#), [\[4\]](#)
- [2019 blog post introducing tract](#)