

# Difference of quantization behavior between onednn and cldnn

Felix Kim

# Intro

- A unit test failed since accuracy issue
  - TestName
    - `deconv_scale_actv_quant_u8_eltw_scale_actv_quant_i8.basic/15`
  - Error Message
    - The difference between `ref[i]` and `output_ptr[i]` is 28, which exceeds tolerance, where `ref[i]` evaluates to 25, `output_ptr[i]` evaluates to 53, and tolerance evaluates to 2.09999999046325684.

# Operations Review

# Quantization

- Meaning

- Quantization is the process of mapping input values from a large set (often a continuous set) to output values in a (countable) smaller set, often with a finite number of elements

- Formula  $output = \frac{\text{clamp}(input; input\_low, input\_high) - input\_low}{s} + input\_low$

$$s = \frac{levels - 1}{input\_high - input\_low}$$

$$x_{f32}[:] = scale_x \cdot (x_{int8}[:] - zp_x)$$

- Usecase

- Mixed precision model?

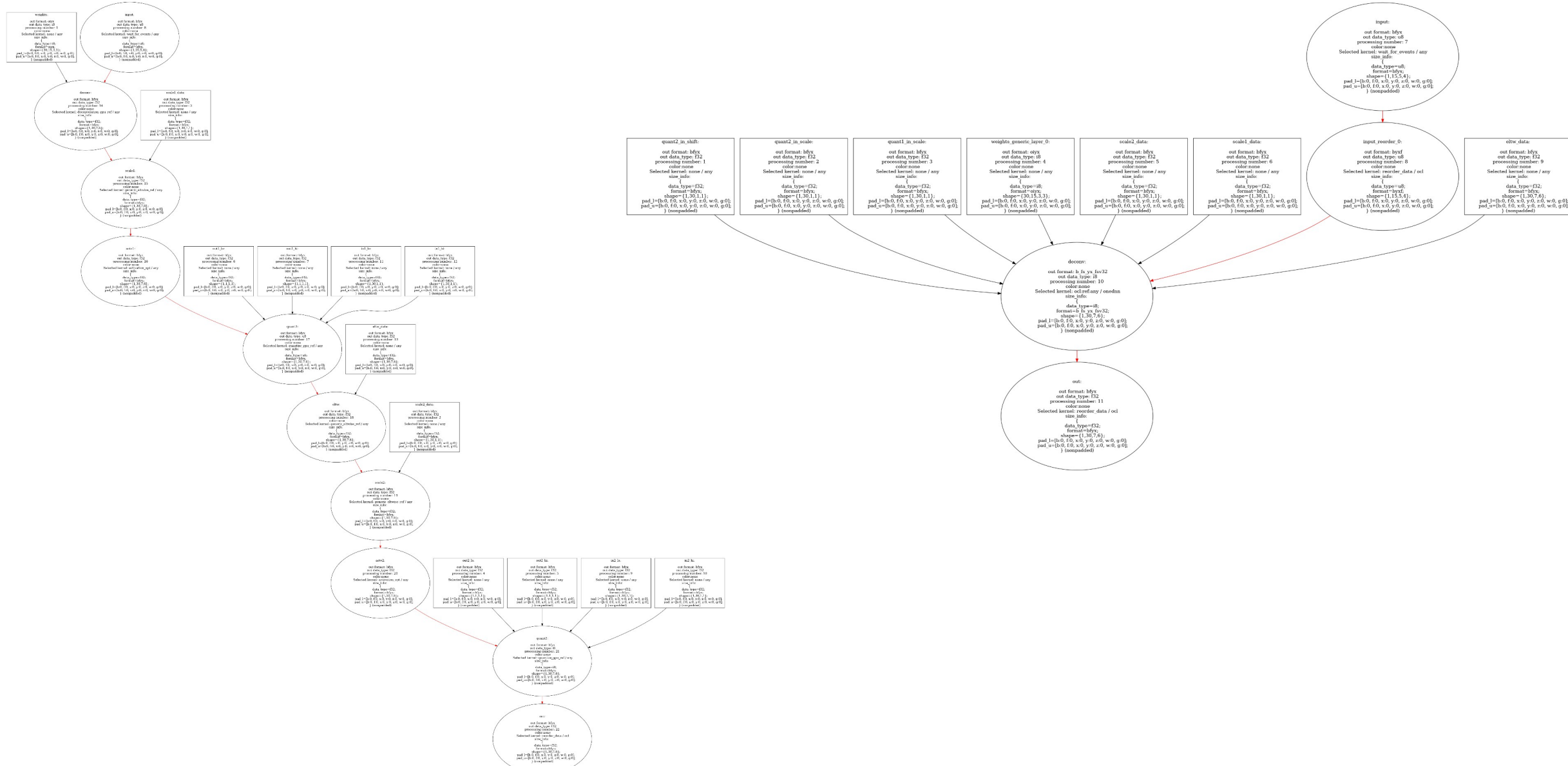
Observations

# Topology

```
input_layout("input", get_input_layout(p)),
data("weights", get_mem(get_weights_layout(p))),
data("scale1_data", get_mem(get_per_channel_layout(p), 1.f / p.kernel.count())),
data("in1_lo", get_mem(get_per_channel_layout(p), 0)),
data("in1_hi", get_mem(get_per_channel_layout(p), 1, max_random)),
data("out1_lo", get_mem(get_single_element_layout(p), 0)),
data("out1_hi", get_mem(get_single_element_layout(p), 255)),
data("eltw_data", get_mem(layout(p.default_type, p.input_format, p.out_shape))),
data("scale2_data", get_mem(get_per_channel_layout(p), 1.f / p.kernel.count())),
data("in2_lo", get_mem(get_per_channel_layout(p), min_random, 0)),
data("in2_hi", get_mem(get_per_channel_layout(p), 1, max_random)),
data("out2_lo", get_mem(get_single_element_layout(p), -127)),
data("out2_hi", get_mem(get_single_element_layout(p), 127)),
```

```
deconvolution("deconv", input_info("input"), { "weights" }, p.groups, p.stride, p.pad),
eltwise("scale1", { input_info("deconv"), input_info("scale1_data") }, eltwise_mode::prod),
activation("actv1", input_info("scale1"), activation_func::relu),
quantize("quant1", input_info("actv1"), input_info("in1_lo"), input_info("in1_hi"),
| | | input_info("out1_lo"), input_info("out1_hi"), 256, data_types::u8),
eltwise("eltw", { input_info("quant1"), input_info("eltw_data") }, eltwise_mode::sum, p.default_type),
eltwise("scale2", { input_info("eltw"), input_info("scale2_data") }, eltwise_mode::prod),
activation("actv2", input_info("scale2"), activation_func::relu),
quantize("quant2", input_info("actv2"), input_info("in2_lo"), input_info("in2_hi"),
| | | input_info("out2_lo"), input_info("out2_hi"), 255, data_types::i8),
reorder("out", input_info("quant2"), p.default_format, data_types::f32)
```

# Graph

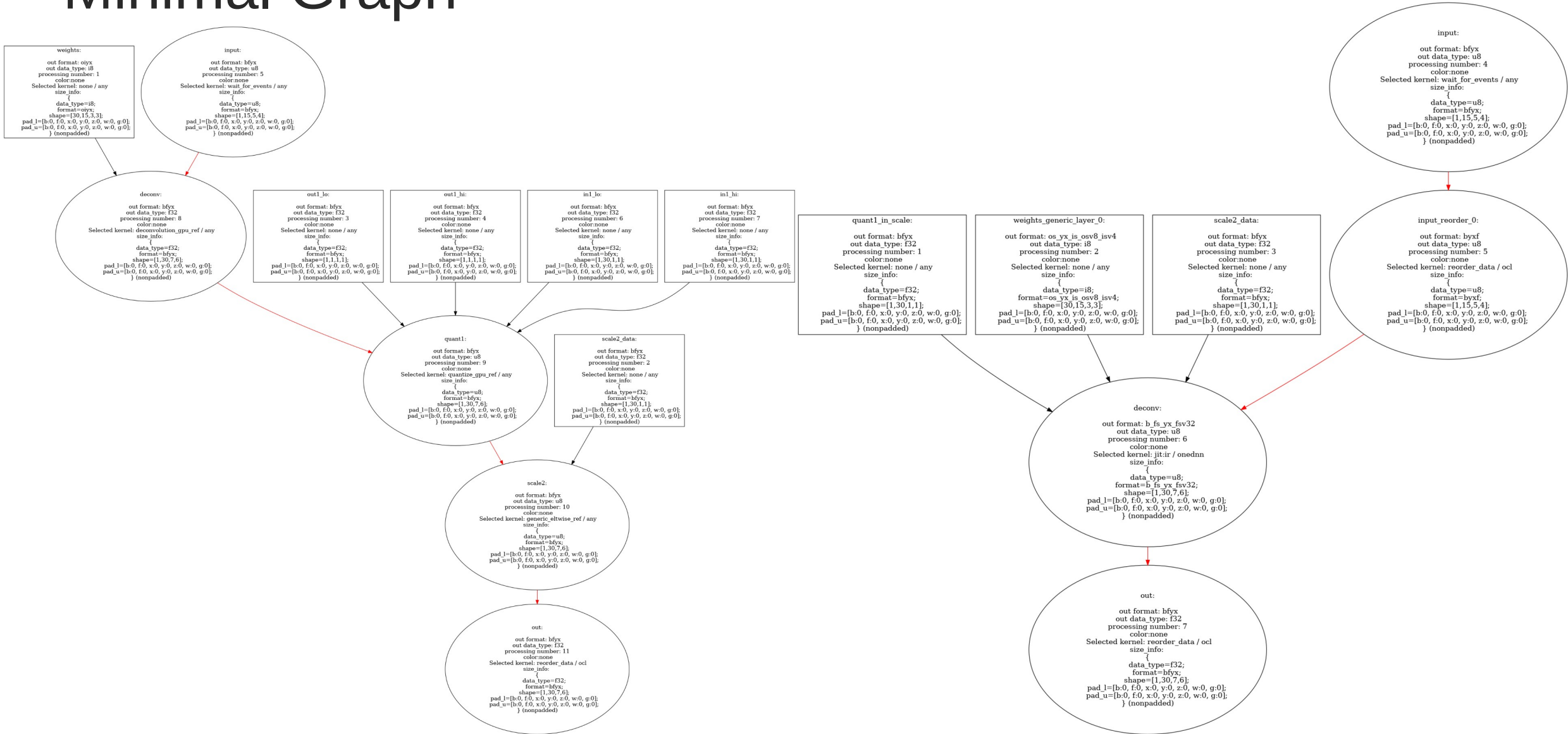


# Minimal Topology

```
input_layout("input", get_input_layout(p)),
data("weights", get_mem(get_weights_layout(p))),
data("in1_lo", get_mem(get_per_channel_layout(p), 0)),
data("in1_hi", get_mem(get_per_channel_layout(p), 1, max_random)),
data("out1_lo", get_mem(get_single_element_layout(p), 0)),
data("out1_hi", get_mem(get_single_element_layout(p), 255)),
data("scale2_data", get_mem(get_per_channel_layout(p), 1.f / p.kernel.count())),
| You, now • Uncommitted changes
deconvolution("deconv", input_info("input"), { "weights" }, p.groups, p.stride, p.pad),
quantize("quant1", input_info("deconv"), input_info("in1_lo"), input_info("in1_hi"),
| | | input_info("out1_lo"), input_info("out1_hi"), 256, data_types::u8),
eltwise("scale2", { input_info("quant1"), input_info("scale2_data") }, eltwise_mode::prod),
reorder("out", input_info("scale2"), p.default_format, data_types::f32)
```



# Minimal Graph



# Output

- Non-optimized
  - 28 28 28 28 28 28 ...
- Optimized
  - 255 255 255 255 255 255 ...
- $255 = \text{max value of u8}$
- $28 = 255 * \text{Scale}$

# Guesses & Conclusion

# Guess: Order of Quantize-Scale is flipped

- Check onednn post-op order
  - DNNL\_VERBOSE=1
  - convert\_dnnl\_verbose.py
- No problem in post-op order
  - eltw\_linear + eltw\_prod
- Interestingly, there is no clip operation

# Root Cause: Clip(Clamp) optimization issue

- For some reason, clip is not added to post-op
  - Has\_clamp turned off in prepare\_quantize\_fusing pass
  - Clamp can be optimized(turned off) when
    - Output datatype = u8/i8
    - Level=256
  - In this case
    - Cldnn use type conversion instead clamp
    - Onednn does not support type conversion post-op

```
bool need_clamp = levels != 256 || out_is_fp;  
bool need_min_clamp = need_clamp; | Ilya Zn  
bool need_max_clamp = need_clamp;
```

# Why other quantize test couldn't catch this?

- Prev primitive = tanh activation
  - Blocked that Quantize being fused
- Level != 256
- Quant is last post-op
- ...

Related codes

# Cldnn round code

```
// Output clamp
if (p->has_clamp) {
    if (p->has_min_clamp && p->has_max_clamp)
        op_decls += "\\n\t" + tmp_var + " = clamp(" + tmp_var + ", " + out_lo + ", " + out_hi + ");";
    else if (p->has_min_clamp)
        op_decls += "\\n\t" + tmp_var + " = max(" + tmp_var + ", " + out_lo + ");";
    else
        op_decls += "\\n\t" + tmp_var + " = min(" + tmp_var + ", " + out_hi + ");";
}

// Output conversion with rounding and saturation
op_decls += "\\n\t" + GetOutputType(vec_size) + " " + out_var + " = " + ConvertToOutputTypeSat(tmp_var, vec_size) + ";";
break;
```



# Fix

- program\_node::init\_onednn\_primitive\_attributes()

```
if (q_param->has_clamp || idx < cldnn_post_ops.size() - 1) {  
    float out_lo = q_param->has_min_clamp ? q_param->out_lo :  
min<float>(out_dt);  
    float out_hi = q_param->has_max_clamp ? q_param->out_hi :  
max<float>(out_dt);  
    post_ops.append_eltwise(1.0f, dnnl::algorithm::eltwise_clip, out_lo,
```

