

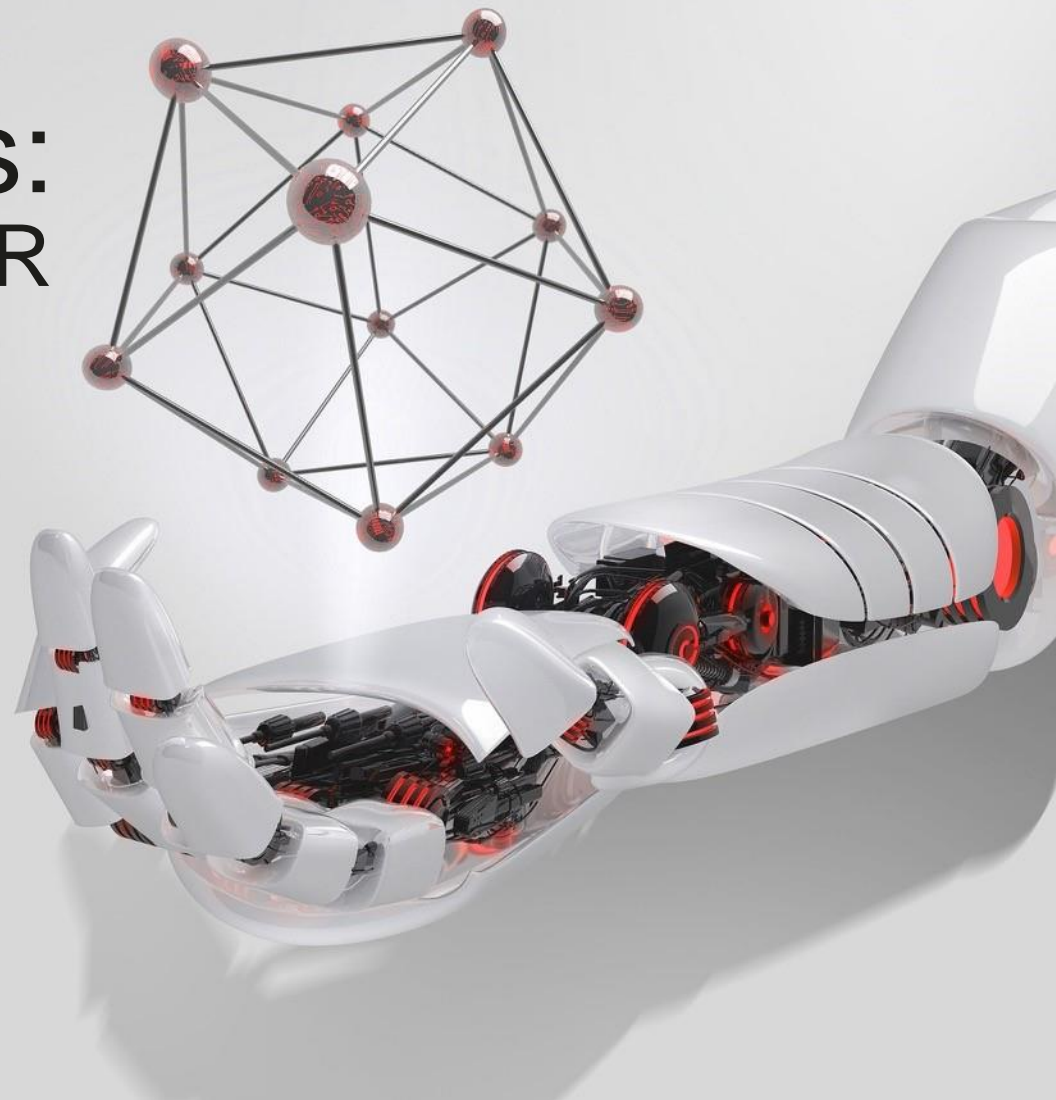
# Transform-dialect schedules: writing MLIR-lowering pipelines in MLIR

**Rolf Morel**

<[rolf.morel@huawei.com](mailto:rolf.morel@huawei.com)>

2024 EuroLLVM, Vienna, April 10-11

Huawei R&D UK, Cambridge  
University of Oxford (DPhil, recently defended)



UNIVERSITY OF  
**OXFORD**

2024 EURO LLVM  
— Vienna Austria —  
DEVELOPERS' MEETING



What we want:

*Declarative & modular* compilers!

Our approach:

- Multi-level rewriting, hence *MLIR*
- Declarative rewriting, hence *schedules*

In this talk:

Schedules for complete lowering & optimization *of* MLIR  
obtained by composing small schedules written *in* MLIR



# What are schedules?



program = algorithm + schedule



# Image processing à la Halide

## algorithm

```
blurx(x,y) = in(x-1,y)
            + in(x,y)
            + in(x+1,y)
out(x,y) = blurx(x,y-1)
          + blurx(x,y)
          + blurx(x,y+1)
```



## program

```
par for out.y0 in 0..out.y.extent/4
  for out.x0 in 0..out.x.extent/4
    alloc blurx[blurx.y.extent][blurx.x.extent]
    for out.yi in 0..4
      let blurx.y.min = 4*out.y0.min + out.yi.min - 1
      for blurx.y in blurx.y.min..blurx.y.max
        for blurx.x0 in blurx.x.min/4..blurx.x.max/4
          vec for blurx.xi in 0..4
            blurx[blurx.y.stride*blurx.y+...] =
              in[in.y.stride*(blurx.y.min+blurx.y)
                +4*blurx.x0+ramp(4)] + ...
          vec for out.xi in 0..4
            out[out.y.stride*(4*(out.y0-out.y0.min)+out.yi)+...] =
              blurx[blurx.y.stride*(out.yi-1-blurx.y.min)
                + out.xi - blurx.x.min] + ...
```



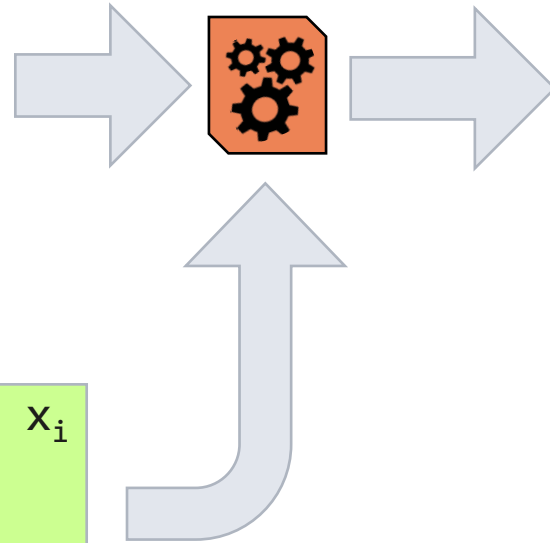
# Image processing à la Halide

**algorithm** (high-level code)

```
blurx(x,y) = in(x-1,y)
             + in(x,y)
             + in(x+1,y)
out(x,y) = blurx(x,y-1)
           + blurx(x,y)
           + blurx(x,y+1)
```

**schedule**

```
blurx: split x by 4 → x0, xi
        vectorize: xi
        store at out.x0
        compute at out.yi
out: split x by 4 → x0, xi
     split y by 4 → y0, yi
     reorder: y0, x0, yi, xi
     parallelize: y0
     vectorize: xi
```



**program** (lowered & optimized code)

```
par for out.y0 in 0..out.y.extent/4
  for out.x0 in 0..out.x.extent/4
    alloc blurx[blurx.y.extent][blurx.x.extent]
    for out.yi in 0..4
      let blurx.y.min = 4*out.y0.min + out.yi.min - 1
      for blurx.y in blurx.y.min..blurx.y.max
        for blurx.x0 in blurx.x.min/4..blurx.x.max/4
          vec for blurx.xi in 0..4
            blurx[blurx.y.stride*blurx.y+...] =
              in[in.y.stride*(blurx.y.min+blurx.y)
                +4*blurx.x0+ramp(4)] + ...
          vec for out.xi in 0..4
            out[out.y.stride*(4*(out.y0-out.y0.min)+out.yi)+...] =
              blurx[blurx.y.stride*(out.yi-1-blurx.y.min)
                + out.xi - blurx.x.min] + ...
```

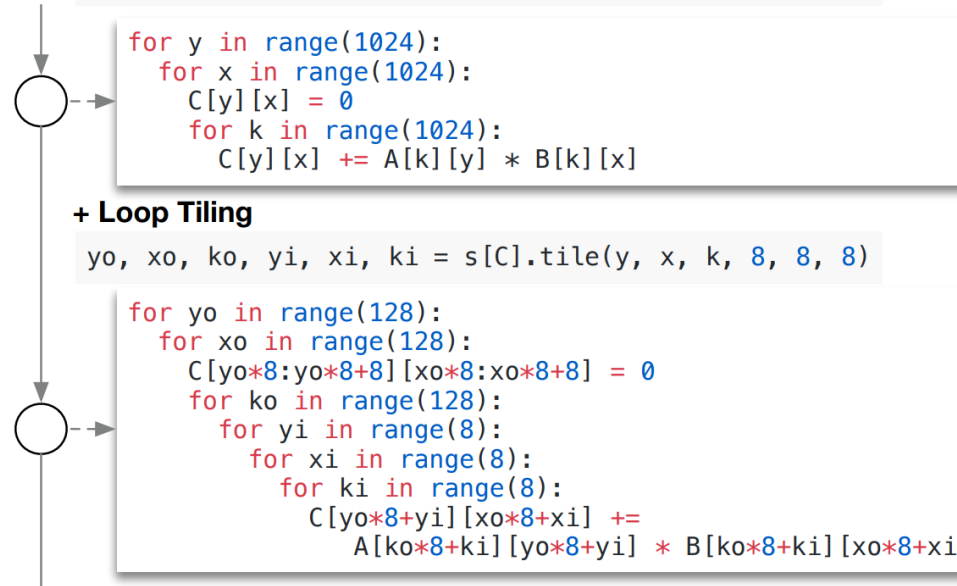
# Many scheduling DSLs

... in support of  
optimizing BLAS / tensor programs:

Tensor Comprehensions  
(Vasilache et al, 2018)

```
A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024), lambda y, x:
              t.sum(A[k, y] * B[k, x], axis=k))
s = t.create_schedule(C.op)
```

TTile  
(Tollenaere et al, 2021)



Fireiron  
(Hagedorn et al, 2020)

Tiramisu  
(Baghdadi et al, 2019)

TVM (Chen et al, 2018)

All share the same *main idea*:  
declarative descriptions of how to transform code



# What is the Transform Dialect?





# MLIR's Transform *meta*-Dialect

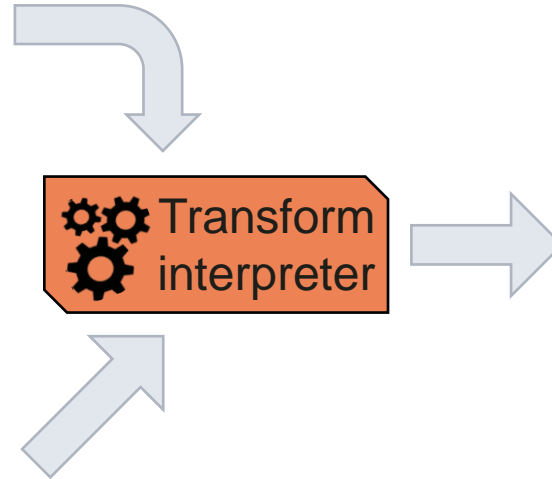
*transform IR* describes transformations on *payload IR*

payload IR:

```
func.func @gemv(%α, %A, %x, %β, %y) {  
  %βy = linalg.generic attrs {iter_types = ["par"]} ... {  
    %βy_elem = arith.mulf %β_elem, %y_elem  
    linalg.yield %βy_elem : f32  
  } -> tensor<?xf32>  
  ...  
  %αAx_plus_βy = linalg.generic  
    { iter_types = ["par", "reduction"]} ... {  
    ...  
  } -> tensor<?x?xf32>  
  return %αAx_plus_βy : tensor<?x?xf32>  
}
```

transform IR:

```
transform.sequence ... {  
  ^bb1(%arg: !transform.any_op):  
    %elemwise = transform.structured.match  
      attrs { iter_types = ["par"] } %arg  
    transform.structured.tile_using_for %elemwise [4]  
}
```



transformed IR:

```
func.func @gemv(%α, %A, %x, %β, %y) {  
  %dim = tensor.dim %y, %c0  
  %βy = scf.for %i = %c0 to %dim step %c4 ... {  
    %ex_slice = tensor.extract_slice ...  
    ... = linalg.generic attrs {iter_types = ["par"]} {  
      %βy_slice_elem = arith.mulf %β_elem, %y_slice_elem  
      linalg.yield %βy_slice_elem : f32  
    } -> tensor<?xf32>  
    %in_slice = tensor.insert_slice ...  
    scf.yield %in_slice  
  }  
  ...  
  %αAx_plus_βy = linalg.generic  
    { iter_types = ["par", "reduction"]} ... {  
    ...  
  } -> tensor<?x?xf32>  
  return %αAx_plus_βy : tensor<?x?xf32>  
}
```



# MLIR's Transform *meta*-Dialect

*transform IR* describes transformations on *payload IR*

payload IR *and* transform IR:

```
module {  
  func.func @name (%arg0, ...) {  
    %0 = some.op(...)  
    %1 = some.other.op(%0, %arg0)  
    ...  
  }  
}  
  
transform.sequence failures(...) {  
  ops for matching ops in payload  
  ...  
  ops for optimizing matched ops  
  AND/OR  
  ops for lowering matched ops  
  ...  
}
```



transformed IR:

```
module {  
  func.func @name (%arg0, ...) {  
    %0 = some.op(...)  
    ...  
    %1 = scf.for %i = %c0 to %arg0 step %k {  
      %2 = some.other.op(%0, %i)  
      ...  
      scf.yield  
    }  
    ...  
  }  
}
```

Tutorial: Controllable Transformations in MLIR, Alex Zinenko, EuroLLVM 2023

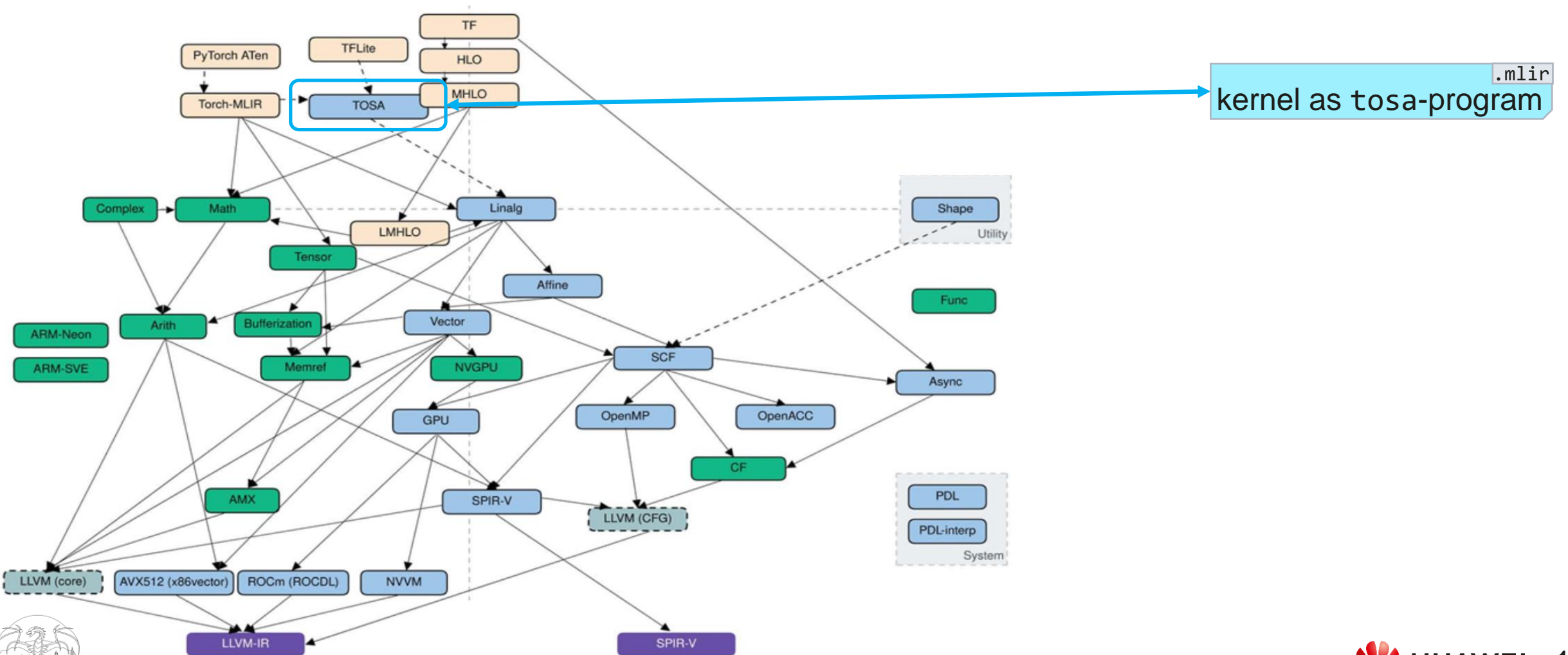
& <https://mlir.llvm.org/docs/Tutorials/transform>, Alex Zinenko



How about combining schedules  
and multi-level rewriting?

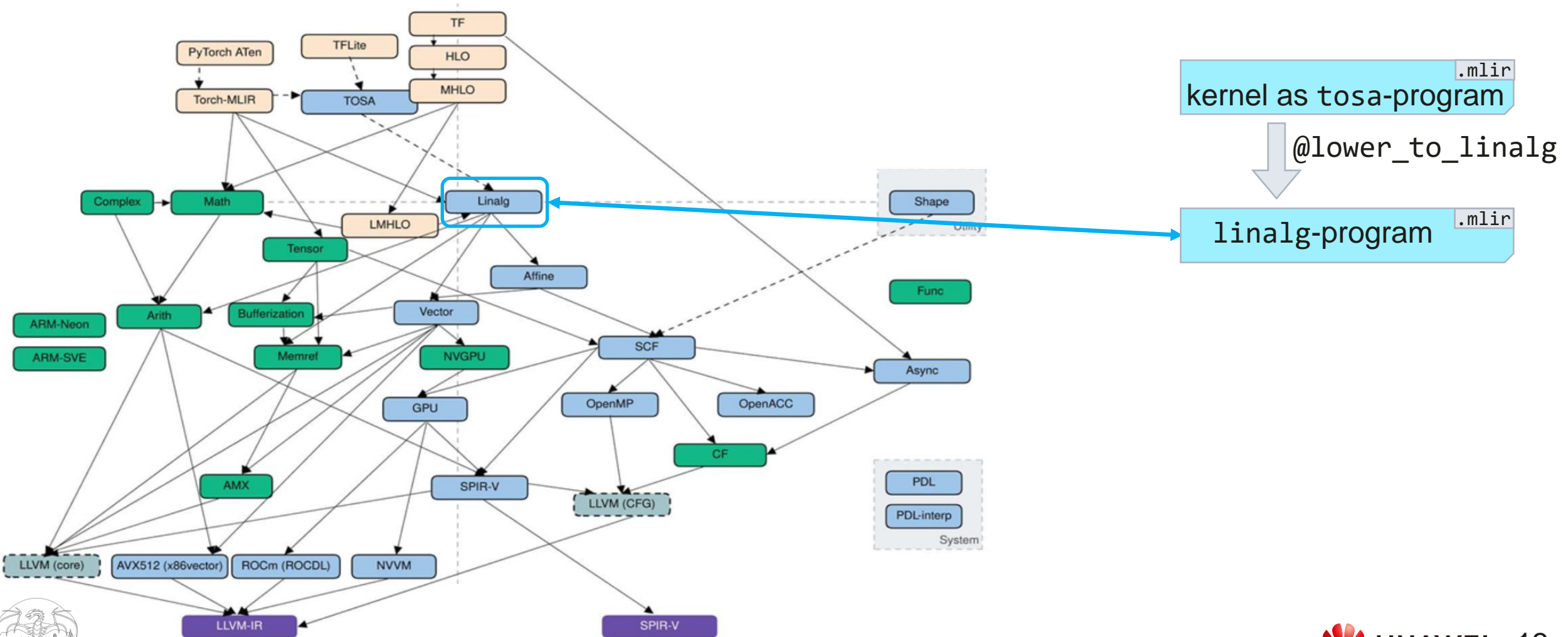


# Progressive lowering/optimization through the dialects



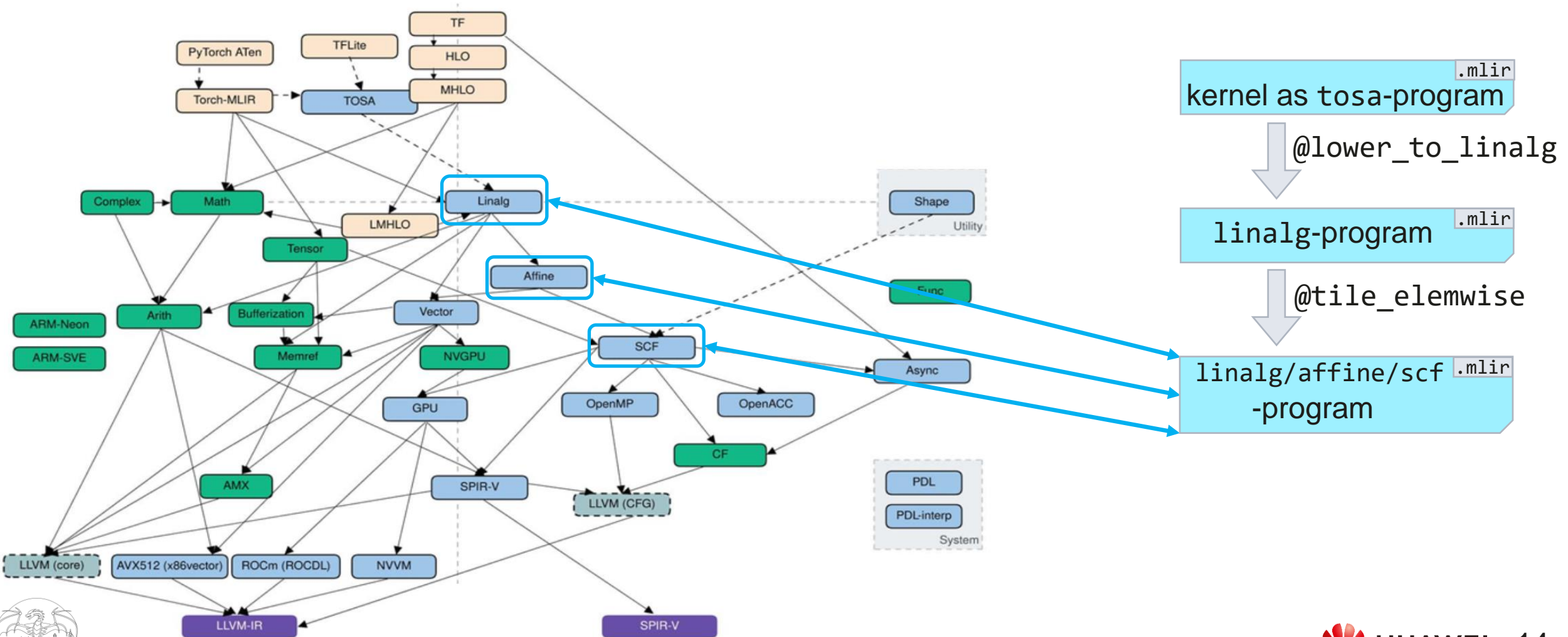
(Quinn Dawkins, 2022)

# Progressive lowering/optimization through the dialects



(Quinn Dawkins, 2022)

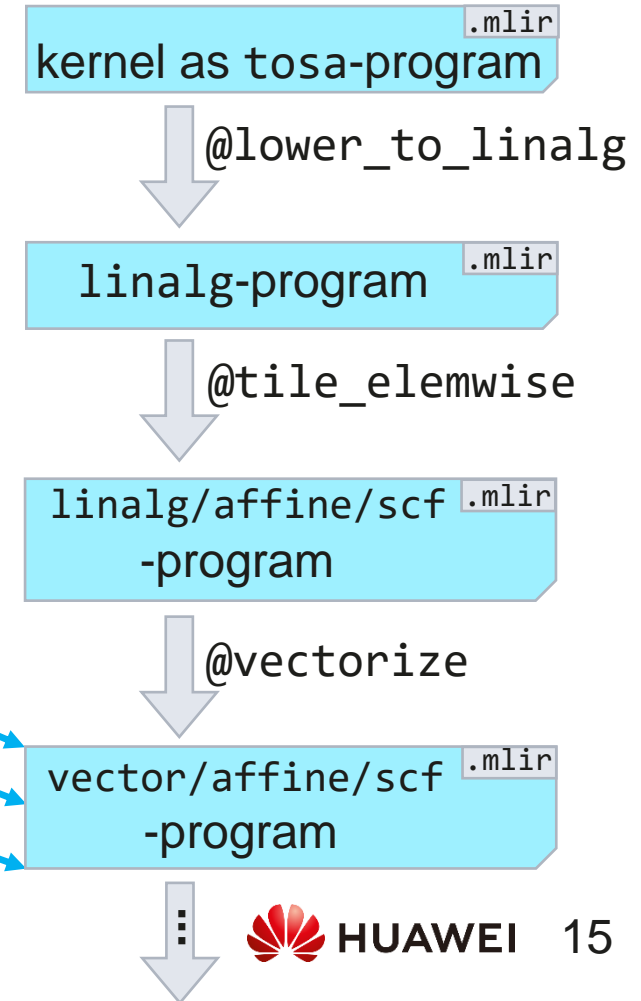
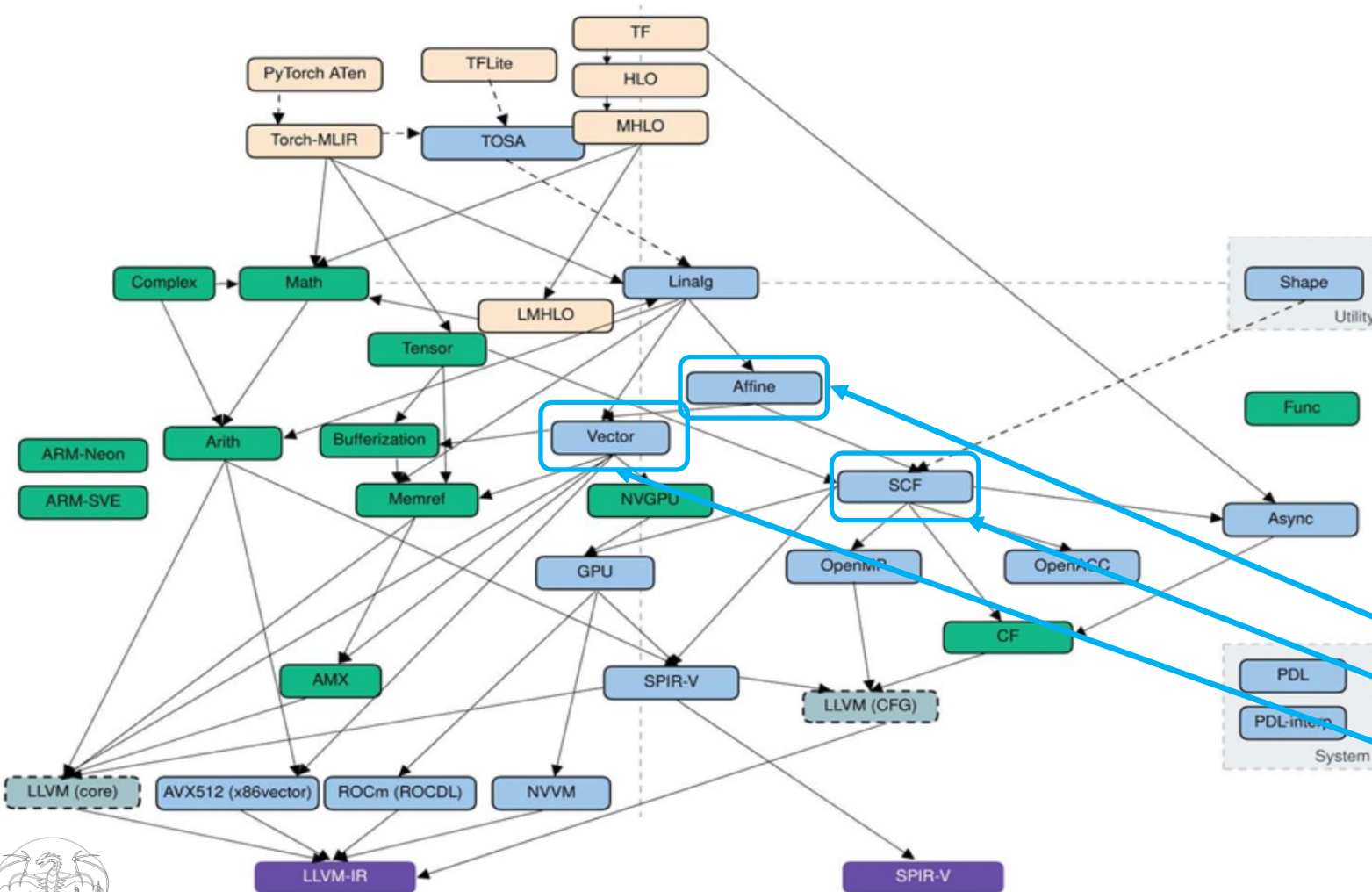
# Progressive lowering/optimization through the dialects



(Quinn Dawkins, 2022)

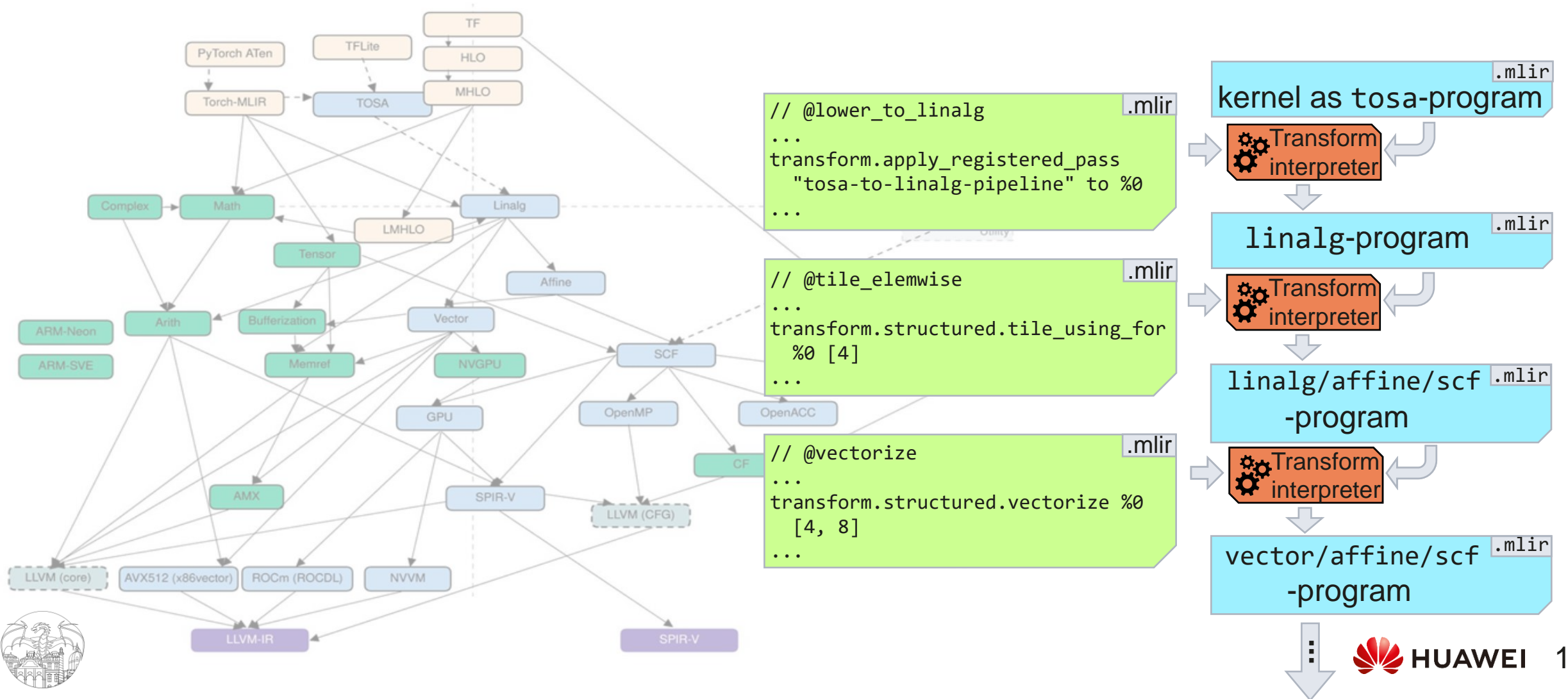


# Progressive lowering/optimization through the dialects



(Quinn Dawkins, 2022)

# Progressive lowering/optimization through the dialects, schedule by schedule

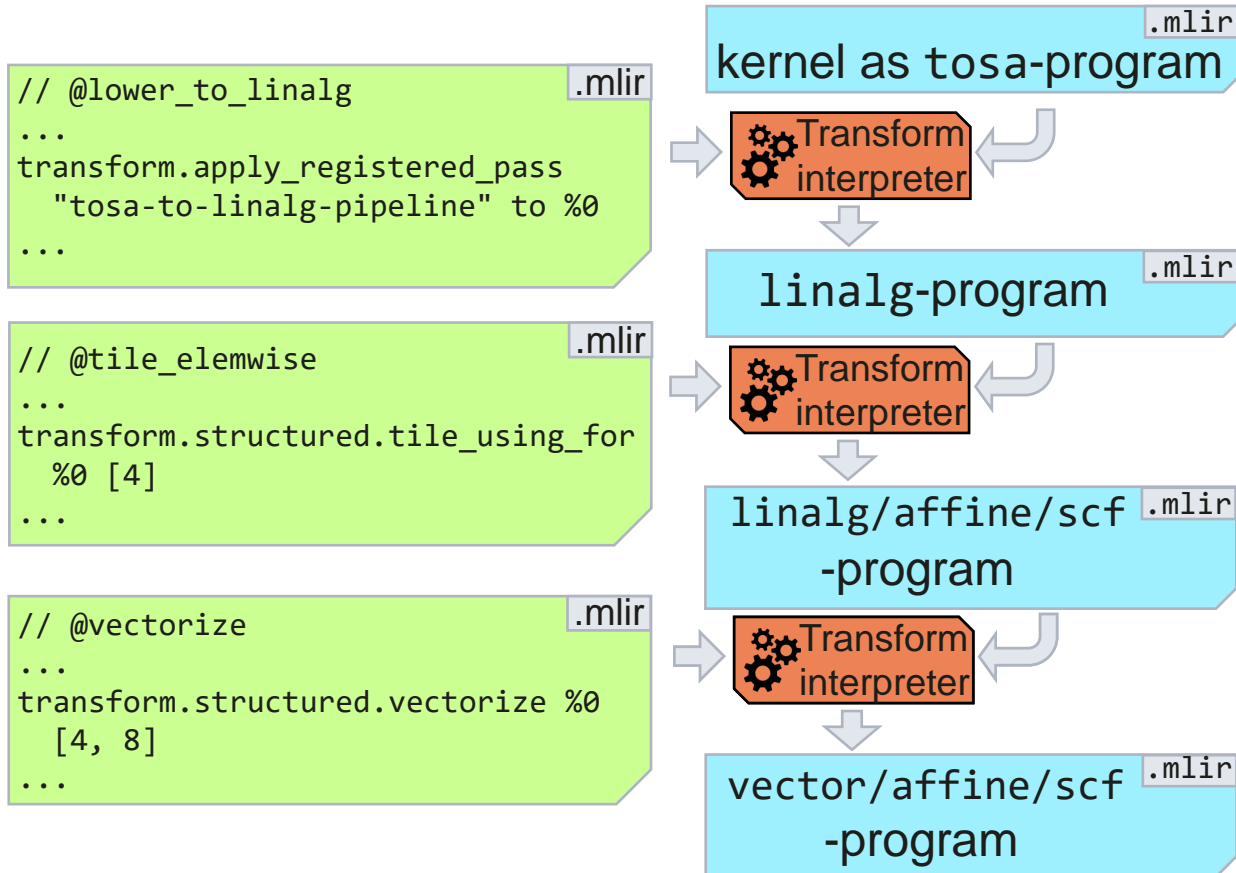




What about composing  
(Transform-dialect) schedules?

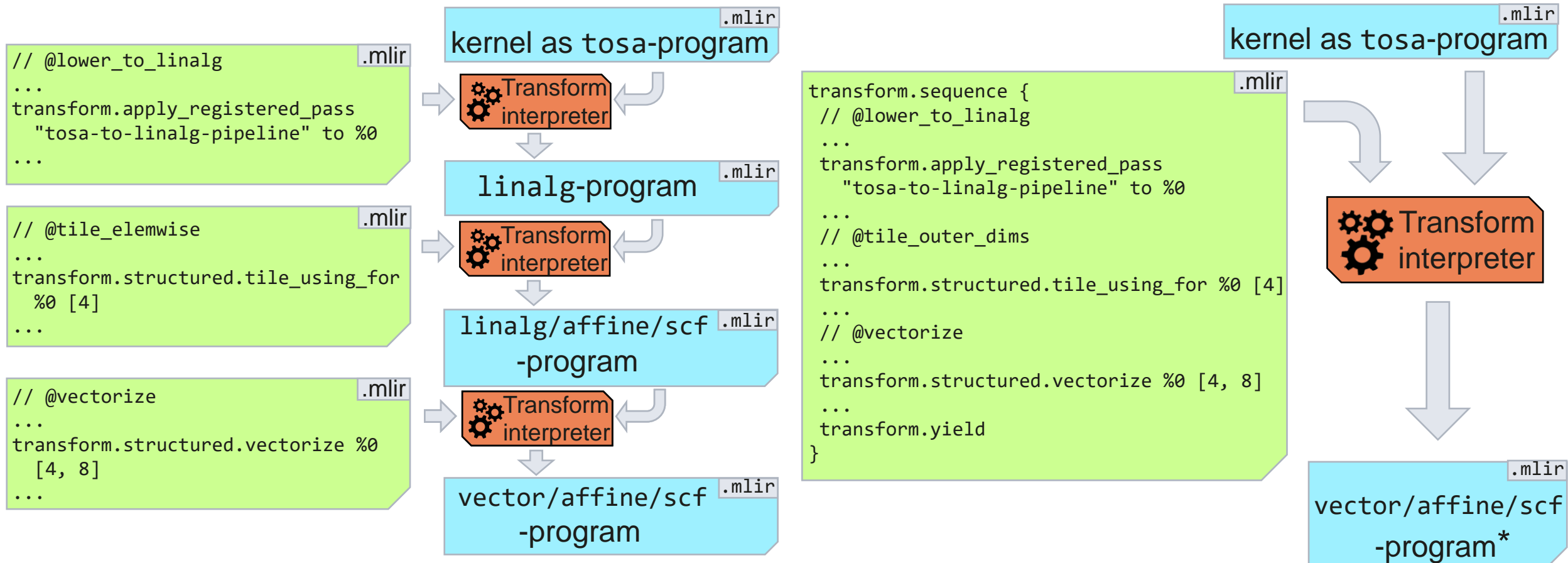


# Monolithic Transform-dialect schedules



\*: might be a different program due to Transform Dialect's backtracking semantics

# Monolithic Transform-dialect schedules



\*: might be a different program due to Transform Dialect's backtracking semantics

# Composable schedules over monolithic schedules

Monolithic sequences are not ideal:

- Need to be programmatically generated
- Structure of the lowering pipeline is lost
- Harder to debug, maintain & reuse



# Composing schedules: Transform Dialect's `named_sequences` and `include`

```
transform.named_sequence @tile_elemwise (%arg0) {  
    %elemwise = transform.structured.match attrs { iter_types = ["par"] } %arg0  
    %tiled, %loop = transform.structured.tile_using_for %elemwise [4]  
    transform.yield %loop  
}  
  
...  
%matched = transform.structured.match %payload  
%tiled_loop = transform.include @tile_elemwise failures(propagate) (%matched)  
...
```



# Composing schedules: main sequence calling other sequences

```
module attributes {transform.with_named_sequence} {  
  transform.named_sequence @lower_to_linalg(%mod) -> !transform.any_op {  
    ...  
    %transformed_mod = ...  
    transform.yield %transformed_mod  
  }  
  ...  
  transform.named_sequence @tile_elemwise(%mod) -> !transform.any_op { ... }  
  transform.named_sequence @vectorize(%mod) -> !transform.any_op { ... }  
  ...  
  transform.named_sequence @__transform_main(%payload) {  
    %mod1 = transform.include @lower_to_linalg failures(propagate) (%payload)  
    %mod2 = transform.include @tile_elemwise failures(propagate) (%mod1)  
    %mod3 = transform.include @vectorize failures(propagate) (%mod2)  
    ...  
  }  
}
```

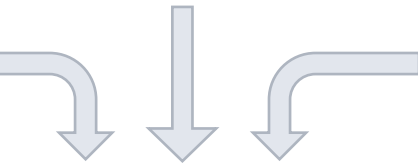
.mlir



# Entire pipelines with reused schedules

```
// kernel1 pipeline .mlir
transform.named_sequence
  @__transform_main(%mod) {
    transform.include @lower_to_linalg ...
    transform.include @tile_elemwise ...
    ...
    transform.include @vectorize ...
    ...
    transform.include @lower_to_llvm ...
  }
```

kernel1 as .mlir  
tosa-program

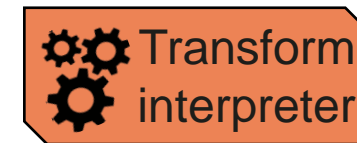
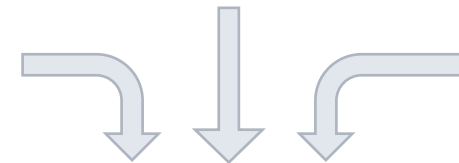


kernel1 as .mlir  
llvm-program

```
transform.named_sequence LIB.mlir
  @lower_to_linalg(%mod) -> !any_op {
    ...
    %transformed_mod = ...
    transform.yield %transformed_mod
  }
transform.named_sequence
  @fuse_reductions(%mod) -> !any_op { ... }
transform.named_sequence
  @tile_elemwise(%mod) -> !any_op { ... }
transform.named_sequence
  @vectorize(%mod) -> !any_op { ... }
...
transform.named_sequence
  @lower_to_llvm(%mod) -> !any_op { ... }
```

```
// kernel2 pipeline .mlir
transform.named_sequence
  @__transform_main(%mod) {
    transform.include @lower_to_linalg ...
    transform.include @fuse_reductions ...
    ...
    transform.include @vectorize ...
    ...
    transform.include @lower_to_llvm ...
  }
```

kernel2 as .mlir  
tosa-program



kernel2 as .mlir  
llvm-program



# Composing schedules ... with glue: CSE and canonicalization

```
transform.named_sequence @__transform_main(%tosa_mod) {  
  ...  
  %mod2 = transform.include @tile_elemwise failures(propagate) (%mod1)  
  %mod2_postcse = transform.apply_cse to %mod2  
  %mod2_postcanon = transform.apply_registered_pass "canonicalize" to %mod2_postcse  
  %mod3_precse = transform.apply_cse to %mod2_postcanon  
  %mod3 = transform.include @vectorize failures(propagate) (%mod3_precse)  
  ...  
}
```

.mlir





So, the Transform Dialect's  
interpreter as a compiler?



# Schedule-based MLIR-compiler

- We use MLIR's Python bindings ...
- ... to generate schedules for each lowering / optimization step
  - ... programmatically so that op attributes, e.g. tile sizes, get set appropriately
  - Mostly upstream Transform ops: less than a dozen are custom
- ... programmatically generate a main sequence for each pipeline
  - Each pipeline lowers from high-level dialects all the way to the LLVM dialect
    - ... using just one schedule ... composed of many small schedules
- ... to delegate running of pipelines fully to the Transform interpreter
  - By invoking a single pass: `-transform-interpreter`



# Schedule reuse in BLAS+-library

We have 19 distinct pipelines (i.e. main sequences) ...  
... which call out to 26 different stepwise schedules

- 15 stepwise schedules are used in 10+ pipelines  
... 4 of these make up a common suffix of all pipelines
- Around 80% of stepwise schedules have at most 7 Transform ops
- Only a couple main sequences have over 11 transform.include ops  
... including the common suffix



# In summary

Schedules allow for declarative descriptions of lowering & optimization

- Transform Dialect allows writing schedules *for MLIR in MLIR*

By writing small composable schedules we can keep our compiler modular

- Small schedules facilitate reuse and maintainability

We can compose schedules for entire pipelines

... and delegate the actual work to the Transform interpreter



# Next steps

- More than just linear pipelines
  - conditional execution: *AlternativesOp* and backtracking upon a match failure
  - One (DAG-shaped) schedule encompassing all pipelines
- Make the stepwise schedules take (e.g. tile size) parameters
  - Would allow for a static `.mlir` library of schedules vs. programmatically generating them
    - Transform ops mainly use `attrs` for parameters (which need to be statically known in MLIR)
- Infer properties of composed schedules, e.g.
  - Inferring overall parameter space for autotuning purposes

