



Dynamic Model in TVM

AWS AI

Presenter: Haichen Shen, Yao Wang
Amazon SageMaker Neo, Deep Engine Science



Models with dynamism

- Control flow (if, loop, etc)
- Dynamic shapes
 - Dynamic inputs: batch size, image size, sequence length, etc.
 - Output shape of some ops are data dependent: arange, nms, etc.
 - Control flow: concatenate within a while loop

Limitation of TVM/graph runtime

- Cannot compile and run dynamic models

Support dynamic model in TVM

- Support Any-dim in typing
- Use shape function to compute the type at runtime
- Virtual machine as a new runtime for Relay
- Dynamic codegen (WIP)
 - Kernel dispatch for a single op
 - Graph dispatch for a (sub-)graph

In collaboration with Jared Roesch, Zhi Chen, Wei Chen



“Any” in Relay typing

Any: represent an unknown dimension at compilation time.

Define a tensor type: `Tensor<(Any, 3, 32, 32), fp32>`

Define type relation:

```
arange: fn(start:fp32, stop:fp32, step:fp32) ->  
Tensor<(Any), fp32>
```

Gradual typing: shape function

- Relax type inference/checking for Any at compilation time

broadcast: `fn(Tensor<(Any, Any), fp32>, Tensor<(1, 8), fp32>) -> Tensor<(Any, 8), fp32>`

Gradual typing: shape function

- Relax type inference/checking for Any at compilation time
- Register a shape function for operator to check the type and compute the output shape

Gradual typing: shape function

- Relax type inference/checking for Any at compilation time
- Register a shape function for operator to check the type and compute the output shape
- Shape function has two modes

(op_attrs, input_tensors, out_ndims) -> out_shape_tensors

- Data dependent

(op_attrs, **input_data**, out_ndims) -> out_shape_tensors

- Data independent

(op_attrs, **input_shapes**, out_ndims) -> out_shape_tensors

Gradual typing: shape function

- Relax type inference/checking for Any at compilation time
- Register a shape function for operator to check the type and compute the output shape
- Shape function has two modes
 - (op_attrs, input_tensors, out_ndims) -> out_shape_tensors**
 - Data dependent
(op_attrs, **input_data**, out_ndims) -> out_shape_tensors
 - Data independent
(op_attrs, **input_shapes**, out_ndims) -> out_shape_tensors
- Why?
 - Fuse data independent shape function together



Shape function example

```
@script
def _concatenate_shape_func(inputs, axis):
    ndim = inputs[0].shape[0]
    out = output_tensor(ndim, "int64")
    for i in const_range(ndim):
        if i != axis:
            out[i] = inputs[0][i]
            for j in const_range(1, len(inputs)):
                assert out[i] == inputs[j][i], "Dims mismatch in the inputs of concatenate."
        else:
            out[i] = int64(0)
            for j in const_range(len(inputs)):
                out[i] += inputs[j][i]
    return out

@_reg.register_shape_func("concatenate", False)
def concatenate_shape_func(attrs, inputs, _):
    axis = get_const_int(attrs.axis)
    return [_concatenate_shape_func(inputs, convert(axis))]
```

Use hybrid script to write shape function

Input shape tensors

Type checking

Data independent

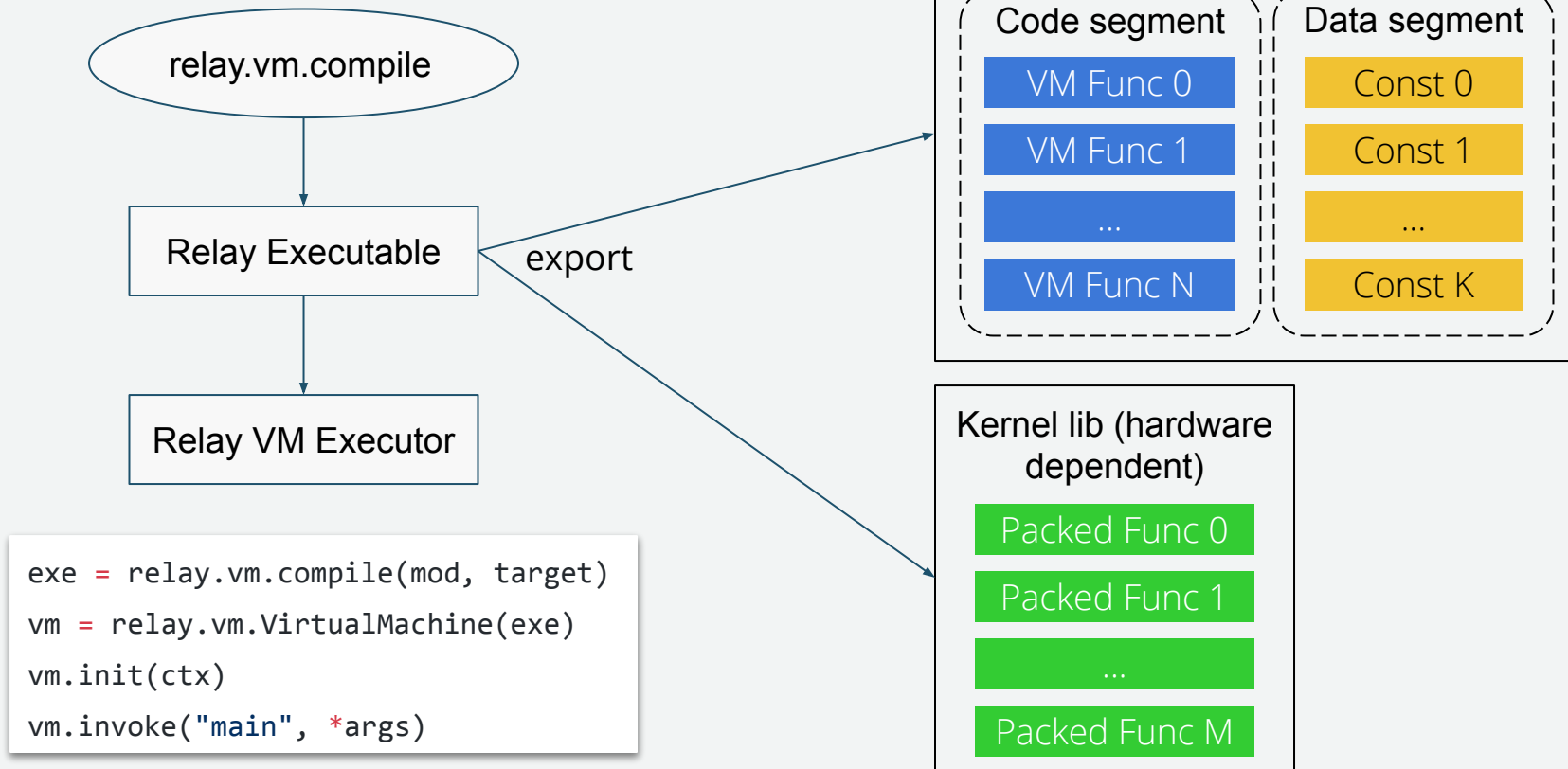
Shape function example

```
@script
def _arange_shape_func(start, stop, step):
    out = output_tensor((1,), "int64")
    out[0] = int64(ceil_div((int64(stop[0]) - int64(start[0])), int64(step[0])))
    return out
```

```
@_reg.register_shape_func("arange", True) Data dependent
def arange_shape_func(attrs, inputs, _):
    return [_arange_shape_func(*inputs)]
```



Relay virtual machine



VM bytecode

Instruction	Description
Move	Moves data from one register to another.
Ret	Returns the object in register result to caller's register.
Invoke	Invokes a function at in index.
InvokeClosure	Invokes a Relay closure.
InvokePacked	Invokes a TVM compiled kernel.
AllocStorage	Allocates a storage block.
AllocTensor	Allocates a tensor value of a certain shape.
AllocTensorReg	Allocates a tensor based on a register.
AllocDatatype	Allocates a data type using the entries from a register.
AllocClosure	Allocates a closure with a lowered virtual machine function.
If	Jumps to the true or false offset depending on the condition.
Goto	Unconditionally jumps to an offset.
LoadConst	Loads a constant at an index from the constant pool.



Relay virtual machine

```
def @main(%i: int32) -> int32 {
  @sum_up(%i) /* ty=int32 */
}

def @sum_up(%i1: int32) -> int32 {
  %0 = equal(%i1, 0 /* ty=int32 */) /* ty=bool */;
  if (%0) {
    %i1
  } else {
    %1 = subtract(%i1, 1 /* ty=int32 */) /* ty=int32 */;
    %2 = @sum_up(%1) /* ty=int32 */;
    add(%2, %i1) /* ty=int32 */
  }
}
```



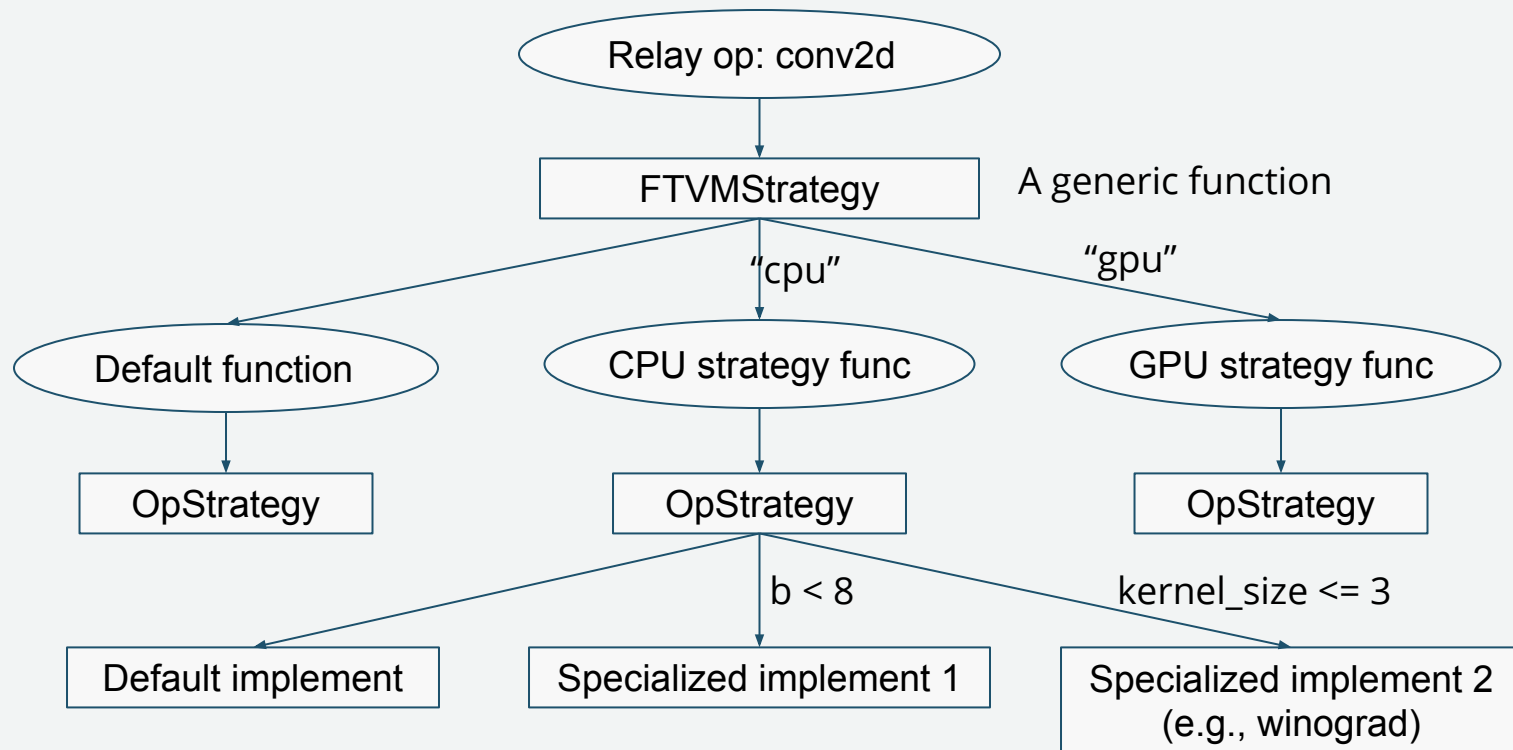
```
sum_up:
  alloc_storage 1 1 64 bool
  alloc_tensor $2 $1 [] uint1
  invoke_packed PackedFunc[0] (in: $0, out: $2)
  load_const $3 1
  if $2 $3 1 2
  goto 9
  alloc_storage 4 4 64 int32
  alloc_tensor $5 $4 [] int32
  invoke_packed PackedFunc[1] (in: $0, out: $5)
  invoke $6 VMFunc[0]($5)
  alloc_storage 7 4 64 int32
  alloc_tensor $8 $7 [] int32
  invoke_packed PackedFunc[2] (in: $6, $0, out: $8)
  move $0 $8
  ret $0

main:
  invoke $1 VMFunc[0]($0)
  ret $1
```

Dynamic codegen: op dispatch (proposal)

- Goal: support codegen for dynamic shape
- Challenges
 - Single kernel performs poor across different shapes
 - Different templates for the same op
 - TVM compute and schedule are coupled together

Dynamic codegen: kernel dispatch (proposal)



Data structure

```
class SpecializedConditionNode : public Node {  
    Array<Expr> conditions;  
};
```

```
class OpImplementNode : public relay::ExprNode {  
    FTVMCompute fcompute;  
    FTVMSchedule fschedule;  
    SpecializedCondition condition; // optional  
};
```

```
class OpStrategyNode : public relay::ExprNode {  
    OpImplement default_implement;  
    Array<OpImplement> specialized_implements;  
};
```

```
class OpStrategy : public relay::Expr {  
    void RegisterDefaultImplement(FTVMCompute fcompute, FTVMSchedule fschedule, bool allow_override=false);  
    void RegisterSpecializedImplement(FTVMCompute fcompute, FTVMSchedule fschedule,  
                                      SpecializedCondition condition);  
};
```



How to register a strategy?

```
@conv2d_strategy.register("cpu")
def conv2d_strategy_cpu(attrs, inputs, out_type, target):
    strategy = OpStrategy()
    layout = attrs.data_layout
    if layout == "NCHW":
        oc, ic, kh, kw = inputs[1].shape
        strategy.register_specialized_implement(wrap_compute_conv2d(topi.x86.conv2d_winograd),
                                                topi.x86.conv2d_winograd,
                                                [kh <= 3, kw <= 3])
        strategy.register_default_implement(wrap_compute_conv2d(topi.x86.conv2d_nchw),
                                            topi.x86.schedule_conv2d_nchw)
    elif layout == "NHWC":
        strategy.register_default_implement(wrap_compute_conv2d(topi.nn.conv2d_nhwc),
                                            topi.x86.schedule_conv2d_nhwc)
    elif layout == "NCHwc":
        strategy.register_default_implement(wrap_compute_conv2d(topi.nn.conv2d_nchwc),
                                            topi.x86.schedule_conv2d_nchwc)
    else: ...
    return strategy
```

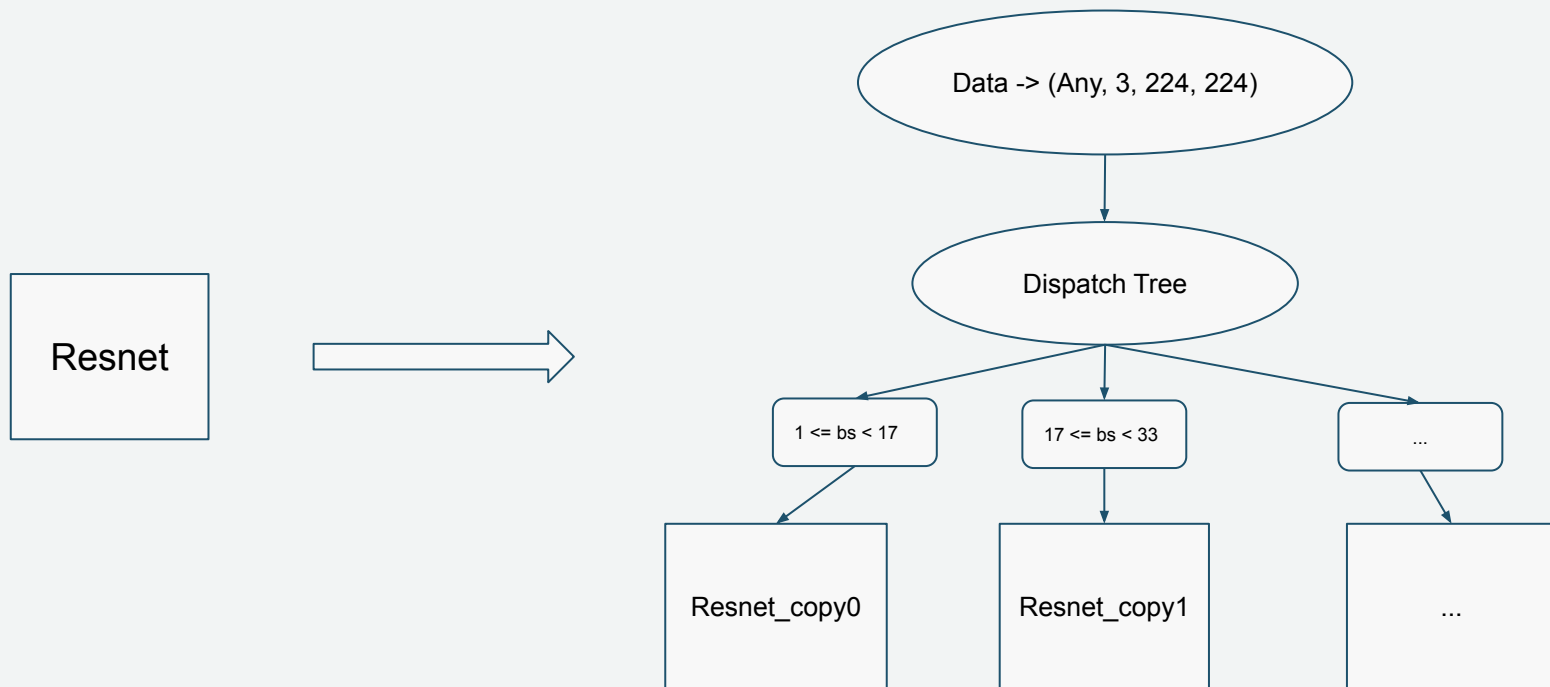


Codegen for OpStrategy

- Each implementation defined will be compiled into a kernel in the module
- Dispatch logic will be compiled into another kernel as well

```
# pseudocode for dispatch kernel
def dispatch_kernel(*args):
    if specialized_condition1:
        specialized_kernel1(*args)
    elif specialized_condition2:
        specialized_kernel2(*args)
    ...
    else:
        default_kernel(*args) # corresponding to default implement
```

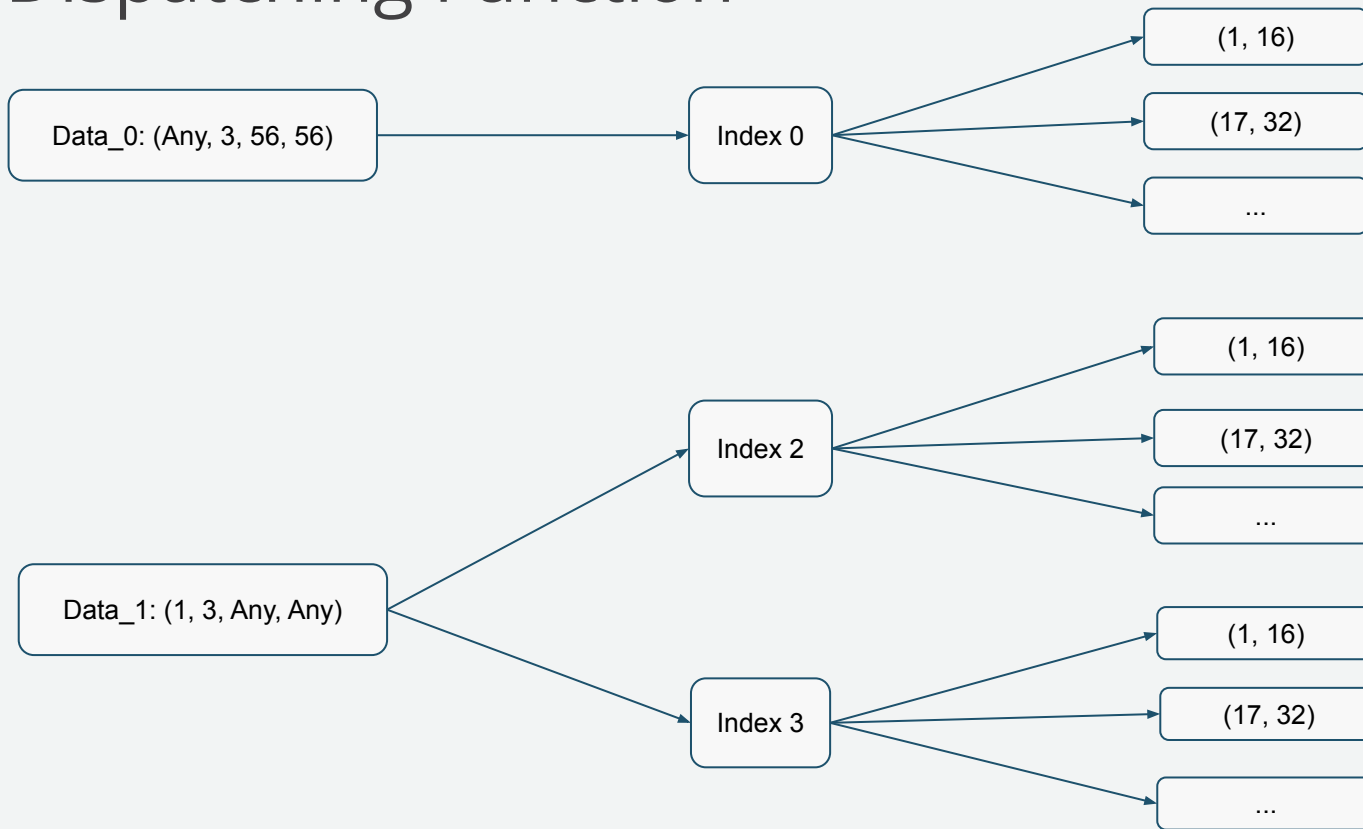
Dispatch a Whole Graph



Why do we need graph dispatcher

1. Minimal overhead: only one dispatching operation is required for each inference.
2. Fit for operator such as conv2d_NCHWc. Graph tuning is well defined for each subgraph.
3. Avoid runtime layout tracking system for operator requires layout transformation to optimize.

Dispatching Function



API Example

```
input_name = "data"
input_shape = [tvm.relay.Any(), 3, 224, 224]
dtype = "float32"
block = get_model('resnet50_v1', pretrained=True)
mod, params = relay.frontend.from_mxnet(block, shape={input_name: input_shape}, dtype=dtype)
tvm.relay.transform.dispatch_global_func(mod, "main", {input_name: input_shape}, tvm.relay.vm.exp_dispatcher)
vmc = relay.backend.vm.VMCompiler()
with tvm.autotvm.apply_graph_best("resnet50_v1_graph_opt.log"):
    vm = vmc.compile(mod, "llvm")

vm.init(ctx)
vm.load_params(params)

data = np.random.uniform(size=(1, 3, 224, 224)).astype("float32")
out = vm.run(data)

data = np.random.uniform(size=(4, 3, 224, 224)).astype("float32")
out = vm.run(data)
```



Acknowledgement



Thank you!

