





**decisions and challenges of writing
an imperative deep learning framework**

Adam Paszke, Sam Gross, Soumith Chintala & the team

Facebook AI Research



Overview of the talk

Computation
Graph Toolkits

Deep Learning
Frameworks

Imperative
Frameworks

JIT Compilation

Declarative

Imperative

Added features
on top of
computation graphs

Design
Challenges
Overheads

Lazy Execution
JIT Fusion
Subgraph caching
differences from AOT

Implementation

Advantages & Disadvantages



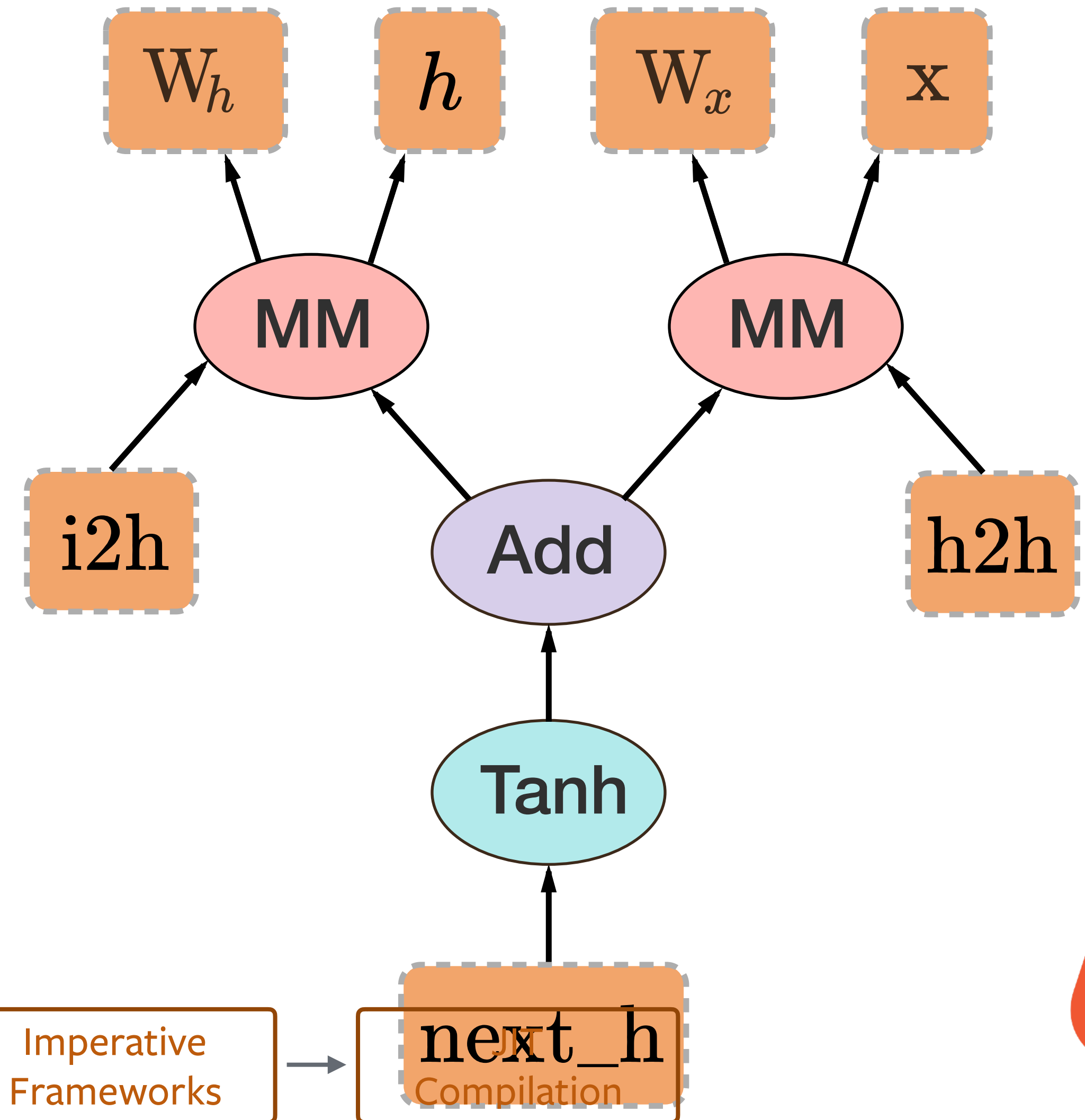
Deep Learning Frameworks



Deep Learning Frameworks

In addition to Computation Graph Toolkits

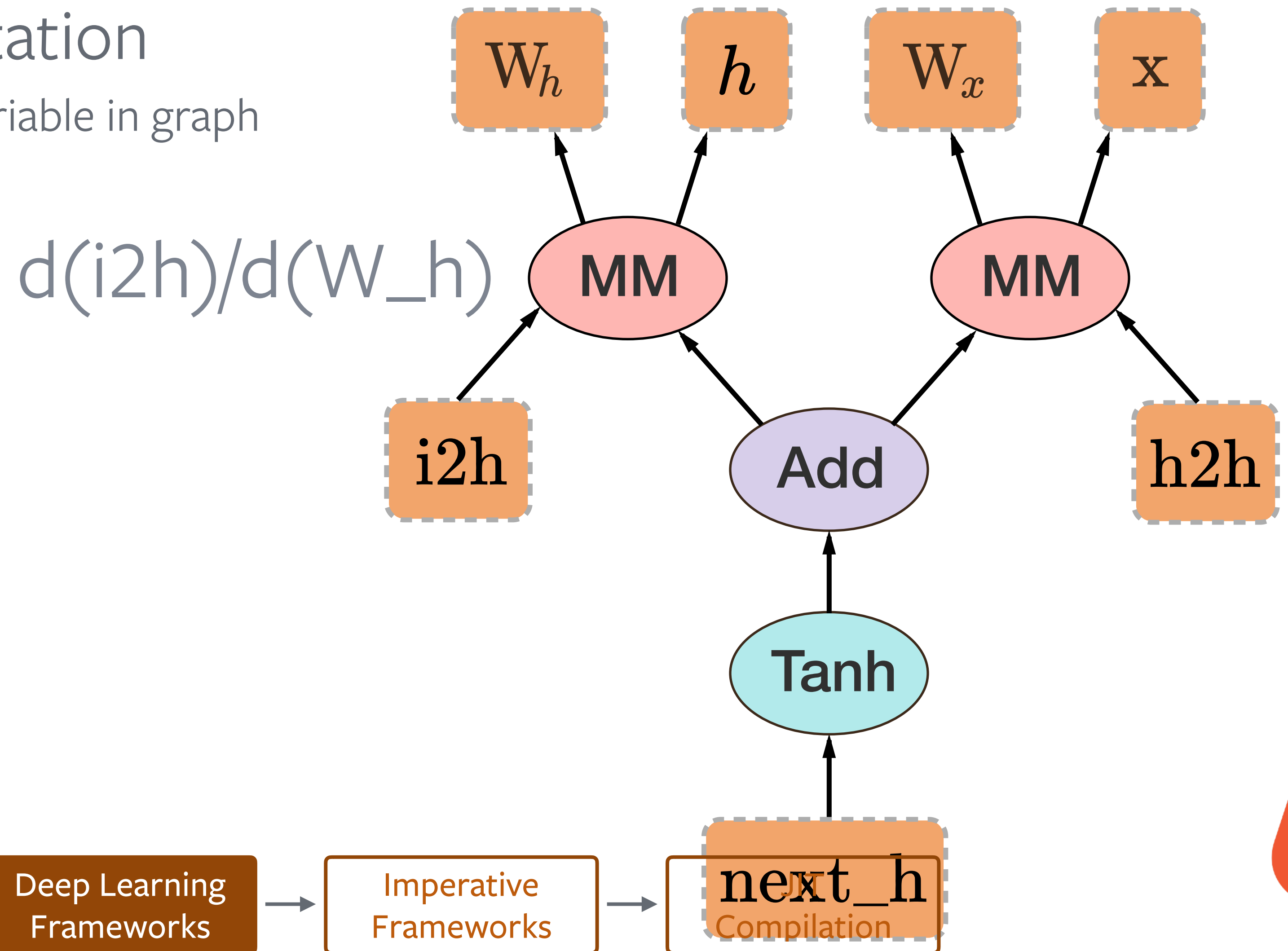
- Provide gradient computation
 - Gradient of one variable w.r.t. any variable in graph



Deep Learning Frameworks

In addition to Computation Graph Toolkits

- Provide gradient computation
 - Gradient of one variable w.r.t. any variable in graph

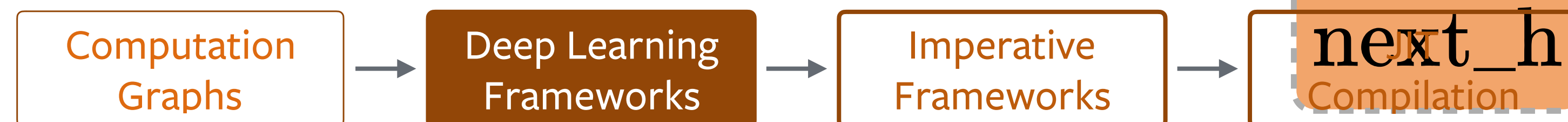
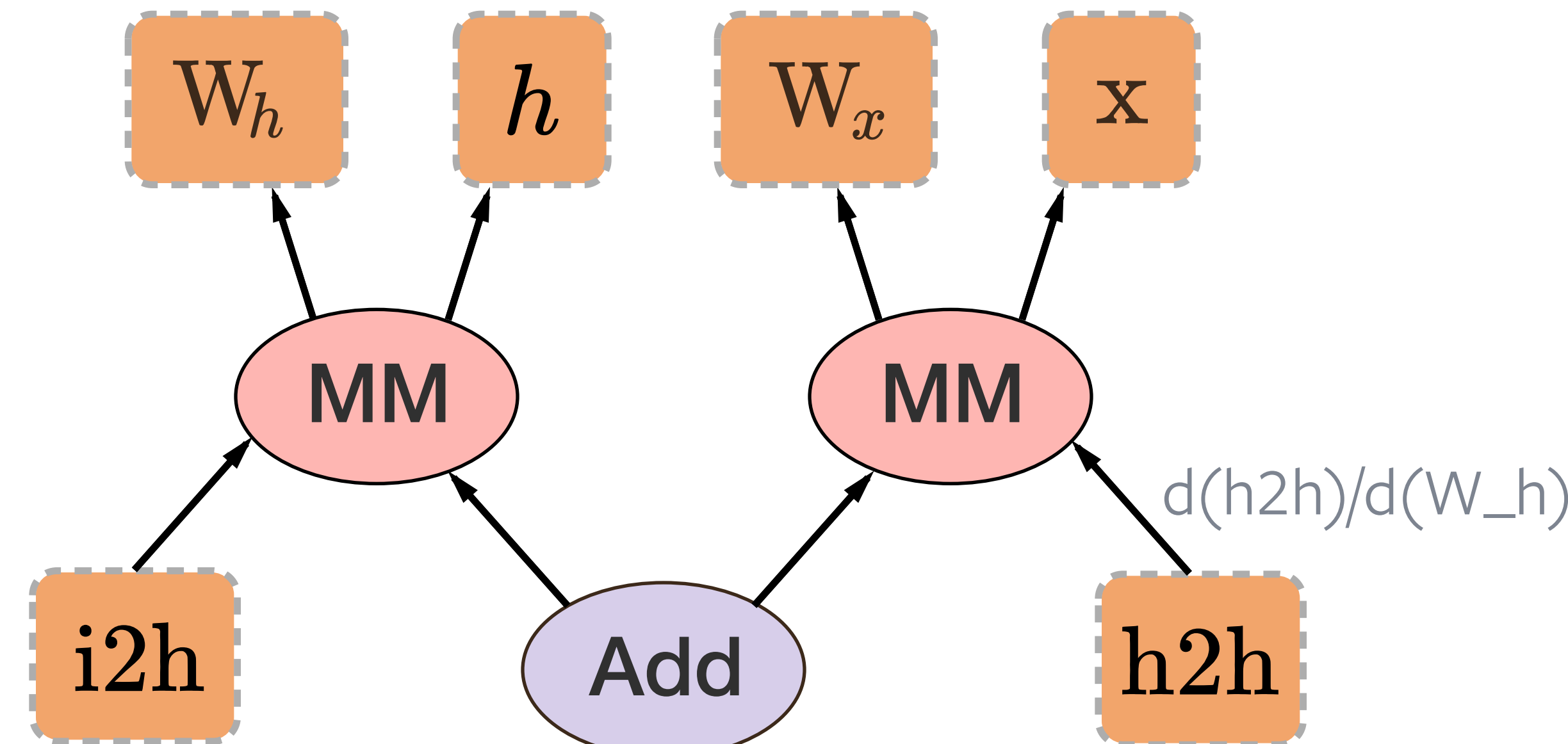


Deep Learning Frameworks

In addition to Computation

Graph Toolkits

- Provide gradient computation
 - Gradient of one variable w.r.t. any variable in graph
- Provide integration with high performance DL libraries like CuDNN



Computation Graph Toolkits



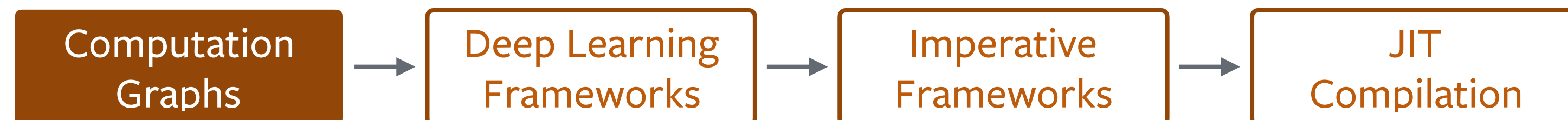
Computation Graph Toolkits

Declarative Toolkits

theano



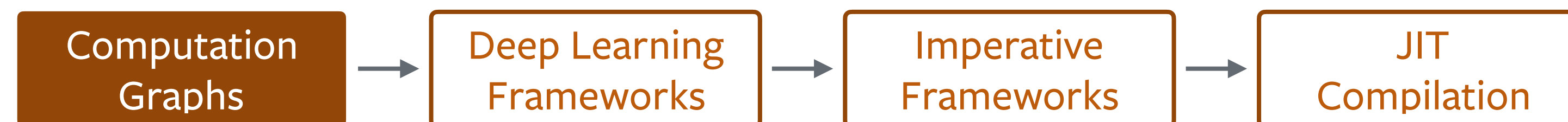
Caffe



Computation Graph Toolkits

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session



Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

```
1 import tensorflow as tf
2 import numpy as np
3
4 trX = np.linspace(-1, 1, 101)
5 trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7 X = tf.placeholder("float")
8 Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27 print(sess.run(w))
```

Computation
Graphs



Deep Learning
Frameworks

Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

A separate turing-complete Virtual Machine

Computation
Graphs

Deep Learning
Frameworks

```
1 import tensorflow as tf
2 import numpy as np
3
4 trX = np.linspace(-1, 1, 101)
5 trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7 X = tf.placeholder("float")
8 Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27 print(sess.run(w))
```


Computation Graph Toolkits

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

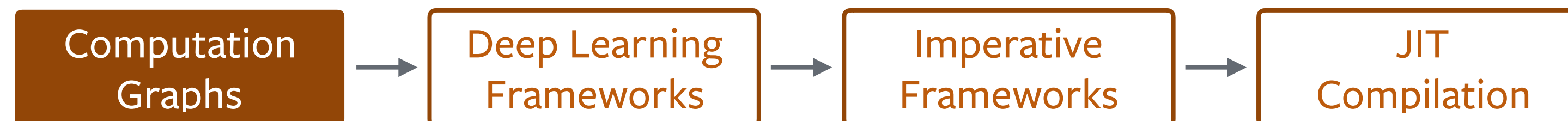
Can handle loops, conditionals
(if, scan, while, etc.)

```
from __future__ import division, print_function
import tensorflow as tf

def fn(previous_output, current_input):
    return previous_output + current_input

elems = tf.Variable([1.0, 2.0, 2.0, 2.0])
elems = tf.identity(elems)
initializer = tf.constant(0.0)
out = tf.scan(fn, elems, initializer=initializer)

with tf.Session() as sess:
    sess.run(tf.initialize_all_variables())
    print(sess.run(out))
```



Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

**Has it's own
execution engine**

Computation
Graphs

Deep Learning
Frameworks

```
1  import tensorflow as tf
2  import numpy as np
3
4  trX = np.linspace(-1, 1, 101)
5  trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7  X = tf.placeholder("float")
8  Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27     print(sess.run(w))
```

Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

Has it's own compiler

- **fuse operations**
- **reuse memory**
- **do optimizations**

Computation
Graphs

Deep Learning
Frameworks

```
1  import tensorflow as tf
2  import numpy as np
3
4  trX = np.linspace(-1, 1, 101)
5  trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7  X = tf.placeholder("float")
8  Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27     print(sess.run(w))
```


Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

Graph can be serialized easily

Computation
Graphs

Deep Learning
Frameworks

```
1 import tensorflow as tf
2 import numpy as np
3
4 trX = np.linspace(-1, 1, 101)
5 trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7 X = tf.placeholder("float")
8 Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27 print(sess.run(w))
```


Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

Own Virtual Machine

Computation
Graphs

Deep Learning
Frameworks

```
1  import tensorflow as tf
2  import numpy as np
3
4  trX = np.linspace(-1, 1, 101)
5  trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7  X = tf.placeholder("float")
8  Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27     print(sess.run(w))
```

Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

Own Virtual Machine
– **separate debugging tools**

Computation
Graphs



Deep Learning
Frameworks

```
1  import tensorflow as tf
2  import numpy as np
3
4  trX = np.linspace(-1, 1, 101)
5  trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7  X = tf.placeholder("float")
8  Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27     print(sess.run(w))
```

Computation Graph

Declarative Toolkits

- Declare a computation
 - with placeholder variables
- Compile it
- Run it in a Session

Own Virtual Machine

- separate debugging tools
- non-linear thinking for user

Computation
Graphs



Deep Learning
Frameworks

```
1 import tensorflow as tf
2 import numpy as np
3
4 trX = np.linspace(-1, 1, 101)
5 trY = 2 * trX + np.random.randn(*trX.shape) * 0.33
6
7 X = tf.placeholder("float")
8 Y = tf.placeholder("float")
9
10 def model(X, w):
11     return tf.multiply(X, w)
12
13 w = tf.Variable(0.0, name="weights")
14 y_model = model(X, w)
15
16 cost = tf.square(Y - y_model)
17
18 train_op = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
19
20 with tf.Session() as sess:
21     tf.global_variables_initializer().run()
22
23     for i in range(100):
24         for (x, y) in zip(trX, trY):
25             sess.run(train_op, feed_dict={X: x, Y: y})
26
27 print(sess.run(w))
```


Imperative Toolkits



Computation Graph Toolkits

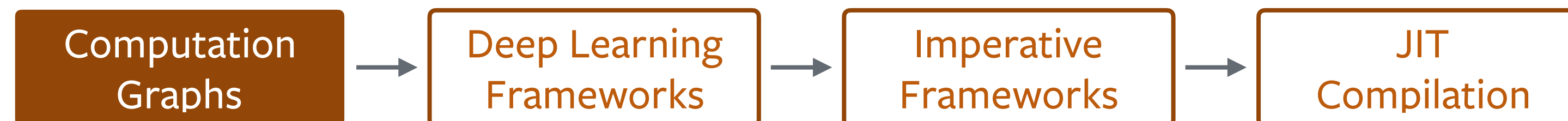
Imperative Toolkits

- Run a computation
- computation is run

PYTORCH

HIPS Autograd

Dynet



Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20  print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20  print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine

```
1 import torch
2 from torch.autograd import Variable
3
4 trX = torch.linspace(-1, 1, 101)
5 trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7 w = Variable(trX.new([0.0]), requires_grad=True)
8
9 for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20 print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy

```
1 import torch
2 from torch.autograd import Variable
3
4 trX = torch.linspace(-1, 1, 101)
5 trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7 w = Variable(trX.new([0.0]), requires_grad=True)
8
9 for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20 print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy

**Oldest debugging method
of all time**

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - Y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20     print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy

Oldest debugging method of all time

```
print(x)
y = foo(x)
print(y)
```

Computation
Graphs

Deep Learning
Frameworks

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13
14         y_model = X * w.expand_as(X)
15         cost = (Y - Y_model) ** 2
16         cost.backward(torch.ones(*cost.size()))
17
18         w.data = w.data + 0.01 * w.grad.data
19
20     print(w)
```

Frameworks

Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy

Oldest debugging method of all time

```
print("hello")  
y = foo(x)  
print("hello2")
```

Computation
Graphs

Deep Learning
Frameworks

```
1  import torch  
2  from torch.autograd import Variable  
3  
4  trX = torch.linspace(-1, 1, 101)  
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33  
6  
7  w = Variable(trX.new([0.0]), requires_grad=True)  
8  
9  for i in range(100):  
10     for (x, y) in zip(trX, trY):  
11         X = Variable(x)  
12         Y = Variable(y)  
13  
14         y_model = X * w.expand_as(X)  
15         cost = (Y - Y_model) ** 2  
16         cost.backward(torch.ones(*cost.size()))  
17  
18         w.data = w.data + 0.01 * w.grad.data  
19  
20     print(w)
```

Frameworks

Compilation

Computation Graph T

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy

**Oldest debugging method
of all time**

```
print("hello")  
y = foo(x)  
print("hello2")
```

Computation
Graphs

Deep Learning
Frameworks

Frameworks

Compilation

```
1  import torch  
2  from torch.autograd import Variable  
3  
4  trX = torch.linspace(-1, 1, 101)  
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33  
6  
7  w = Variable(trX.new([0.0]), requires_grad=True)  
8  
9  for i in range(100):  
10     for (x, y) in zip(trX, trY):  
11         X = Variable(x)  
12         Y = Variable(y)  
13         print(Y)  
14  
15         y_model = X * w.expand_as(X)  
16         cost = (Y - Y_model) ** 2  
17         cost.backward(torch.ones(*cost.size()))  
18  
19         w.data = w.data + 0.01 * w.grad.data  
20  
21     print(w)
```

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- debugging is easy
- Linear program flow
 - Linear thinking for user

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13         print(Y)
14
15         y_model = X * w.expand_as(X)
16         cost = (Y - Y_model) ** 2
17         cost.backward(torch.ones(*cost.size()))
18
19         w.data = w.data + 0.01 * w.grad.data
20
21     print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- Cannot compile program
- Linear program flow
 - Linear thinking for user

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13         print(Y)
14
15         y_model = X * w.expand_as(X)
16         cost = (Y - Y_model) ** 2
17         cost.backward(torch.ones(*cost.size()))
18
19         w.data = w.data + 0.01 * w.grad.data
20
21     print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine
- Cannot compile program
- Cannot optimize
- Linear thinking for user

```
1 import torch
2 from torch.autograd import Variable
3
4 trX = torch.linspace(-1, 1, 101)
5 trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7 w = Variable(trX.new([0.0]), requires_grad=True)
8
9 for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13         print(Y)
14
15     y_model = X * w.expand_as(X)
16     cost = (Y - Y_model) ** 2
17     cost.backward(torch.ones(*cost.size()))
18
19     w.data = w.data + 0.01 * w.grad.data
20
21 print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine

Cannot compile program

Cannot optimize

Cannot do static analysis

```
1 import torch
2 from torch.autograd import Variable
3
4 trX = torch.linspace(-1, 1, 101)
5 trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7 w = Variable(trX.new([0.0]), requires_grad=True)
8
9 for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13         print(Y)
14
15     y_model = X * w.expand_as(X)
16     cost = (Y - Y_model) ** 2
17     cost.backward(torch.ones(*cost.size()))
18
19     w.data = w.data + 0.01 * w.grad.data
20
21 print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

Computation Graph

Imperative Toolkits

- Run a computation
- computation is run!
- no separate execution engine

Cannot compile program

Cannot optimize

Cannot do static analysis
(more on this later)

```
1  import torch
2  from torch.autograd import Variable
3
4  trX = torch.linspace(-1, 1, 101)
5  trY = 2 * trX + torch.randn(*trX.size()) * 0.33
6
7  w = Variable(trX.new([0.0]), requires_grad=True)
8
9  for i in range(100):
10     for (x, y) in zip(trX, trY):
11         X = Variable(x)
12         Y = Variable(y)
13         print(Y)
14
15     y_model = X * w.expand_as(X)
16     cost = (Y - Y_model) ** 2
17     cost.backward(torch.ones(*cost.size()))
18
19     w.data = w.data + 0.01 * w.grad.data
20
21     print(w)
```

Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

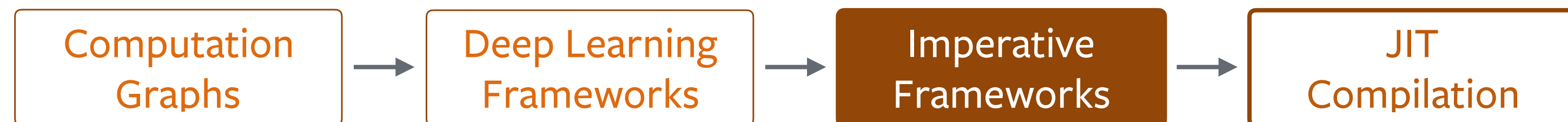
Imperative Frameworks



Imperative Frameworks: PyTorch

Graph is built on the fly

```
from torch.autograd import Variable
```

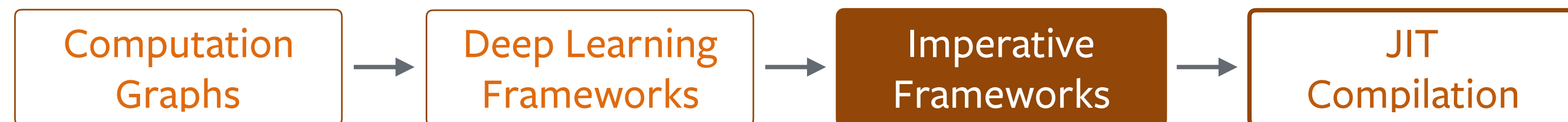
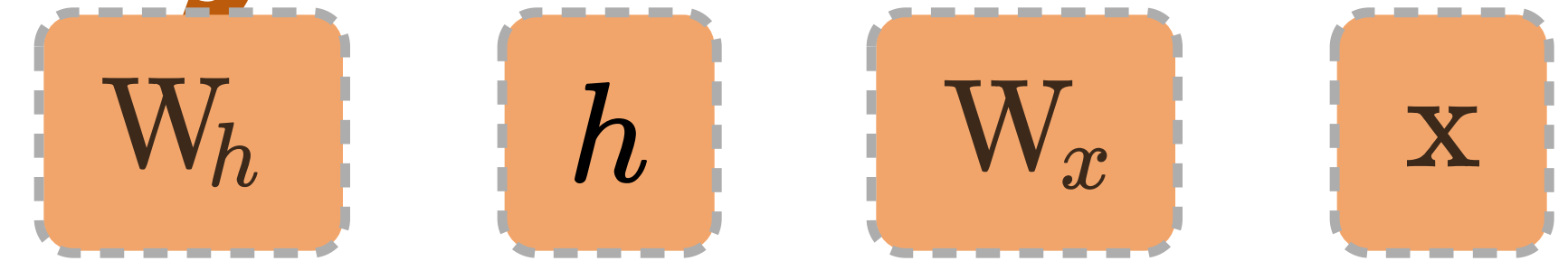


Imperative Frameworks: PyTorch

Graph is built on the fly

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```



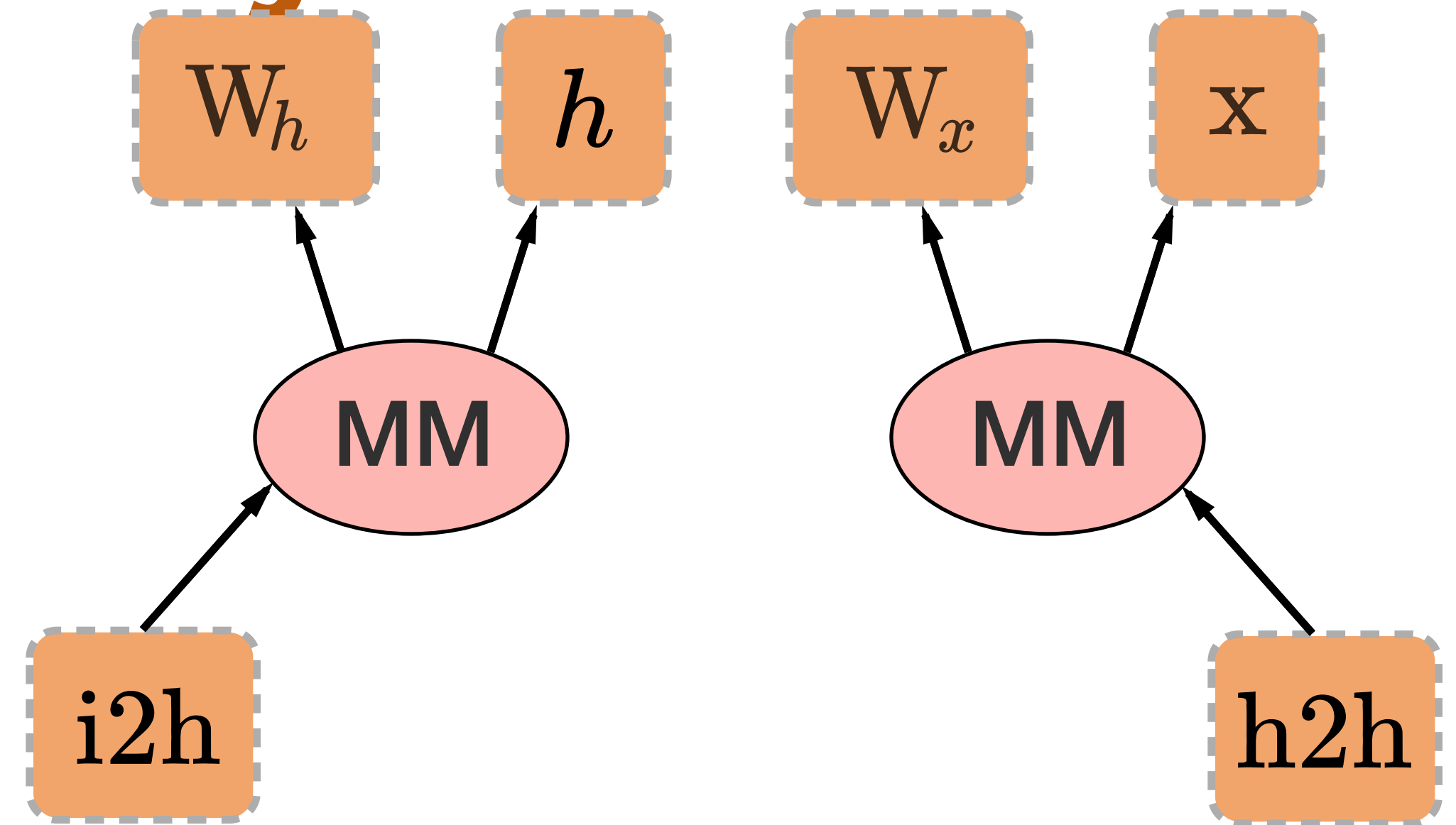
Imperative Frameworks: PyTorch

Graph is built on the fly

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())
```



Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation

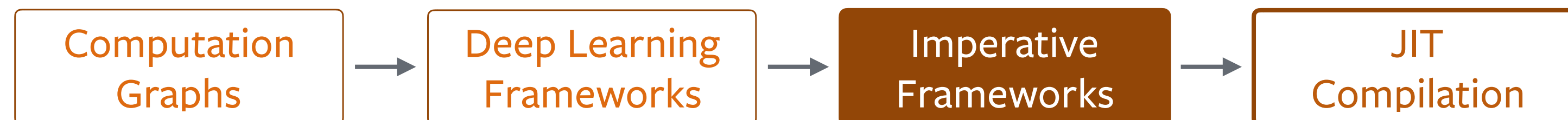
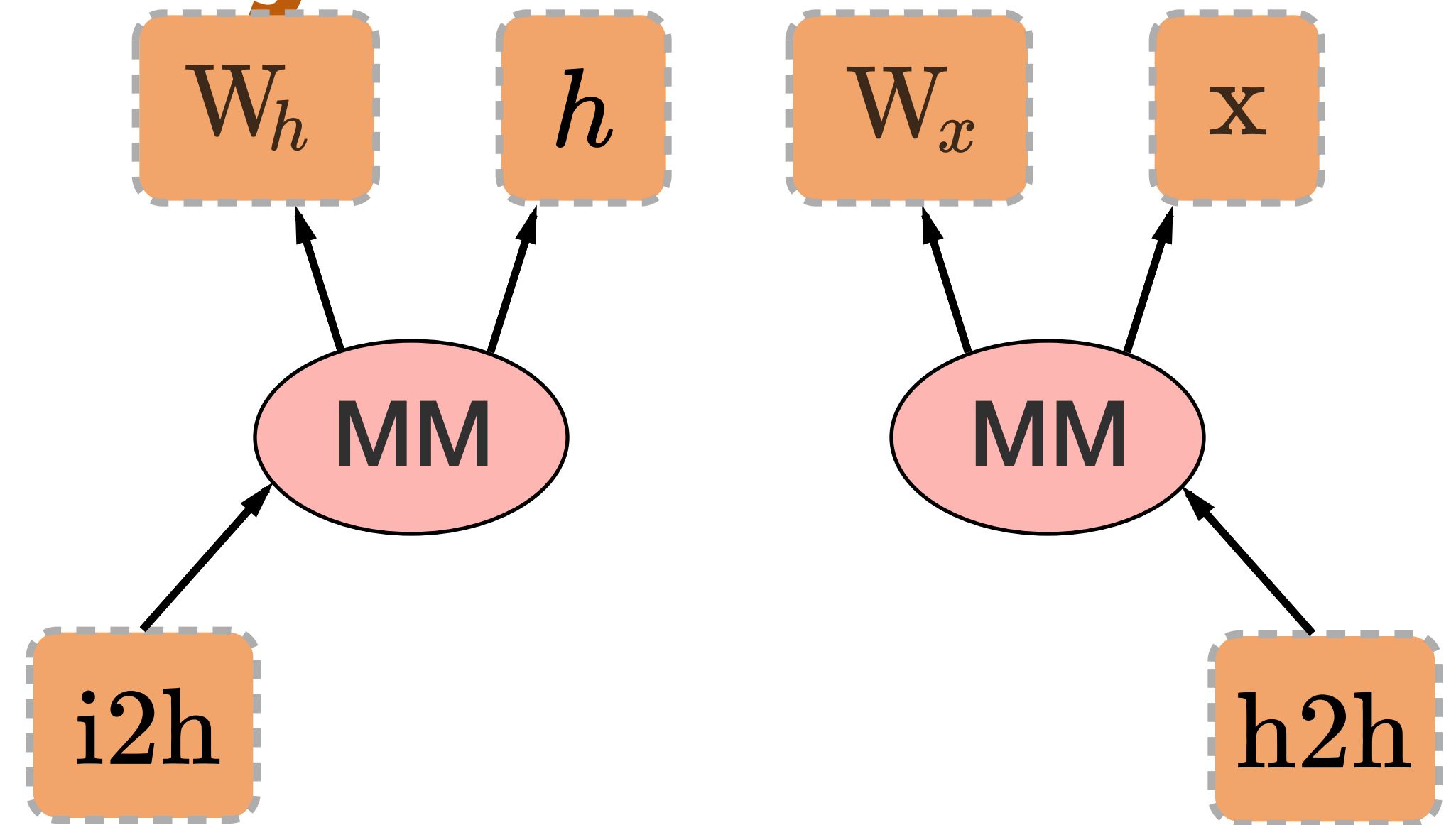


Imperative Frameworks: PyTorch

```
from torch.autograd import Variable

x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))

i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
```

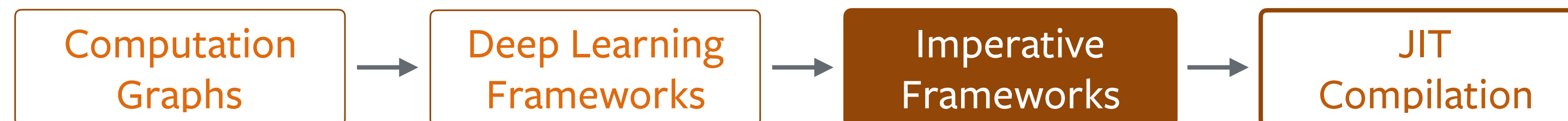
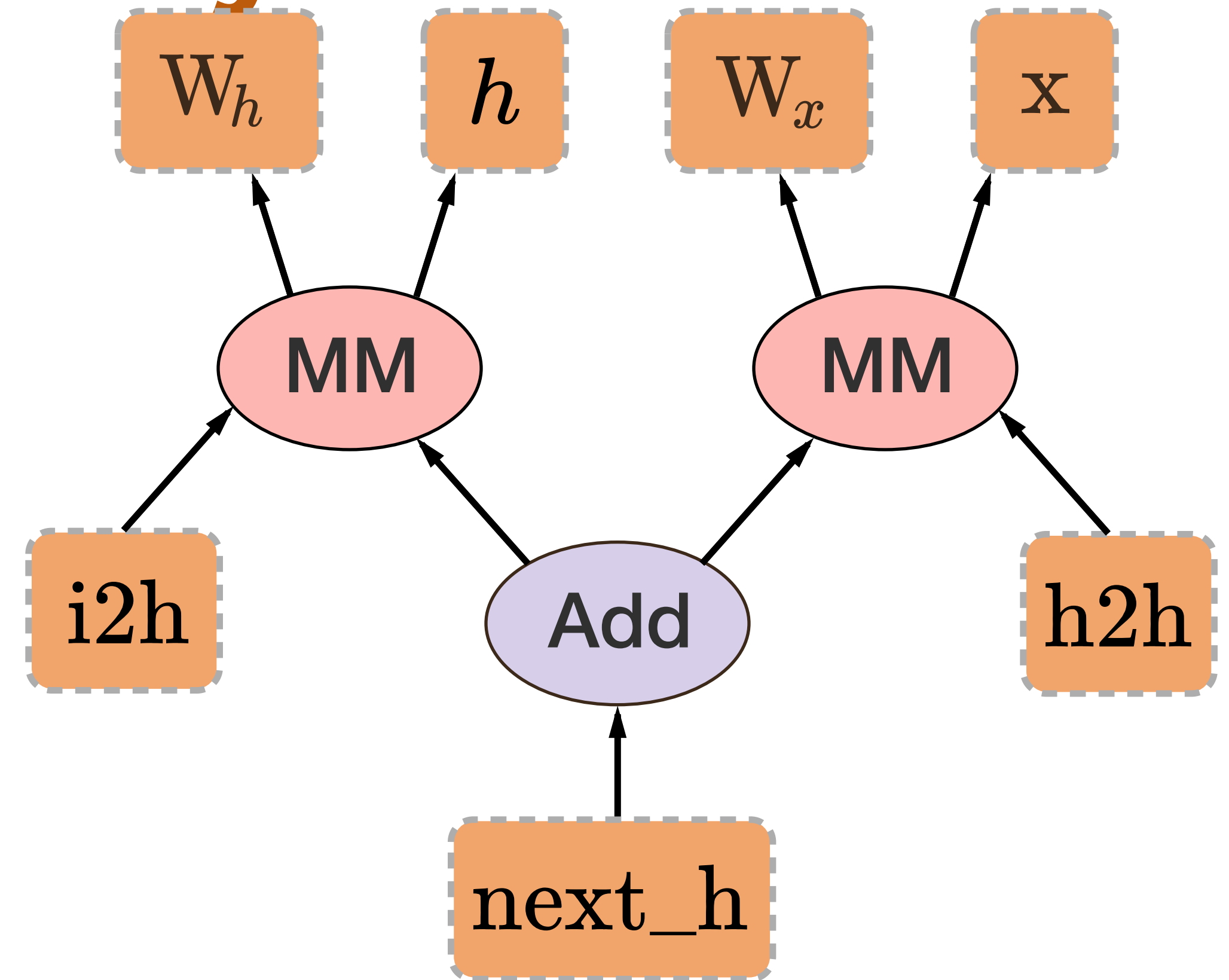


Imperative Frameworks: PyTorch

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h
```

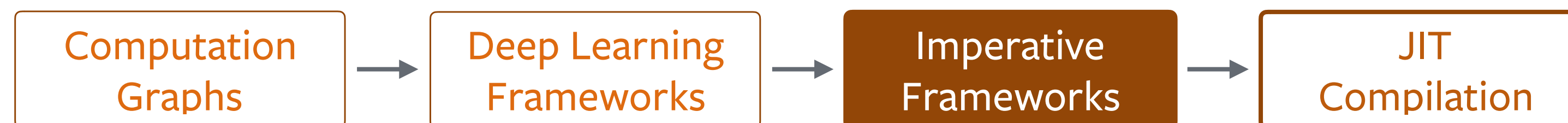
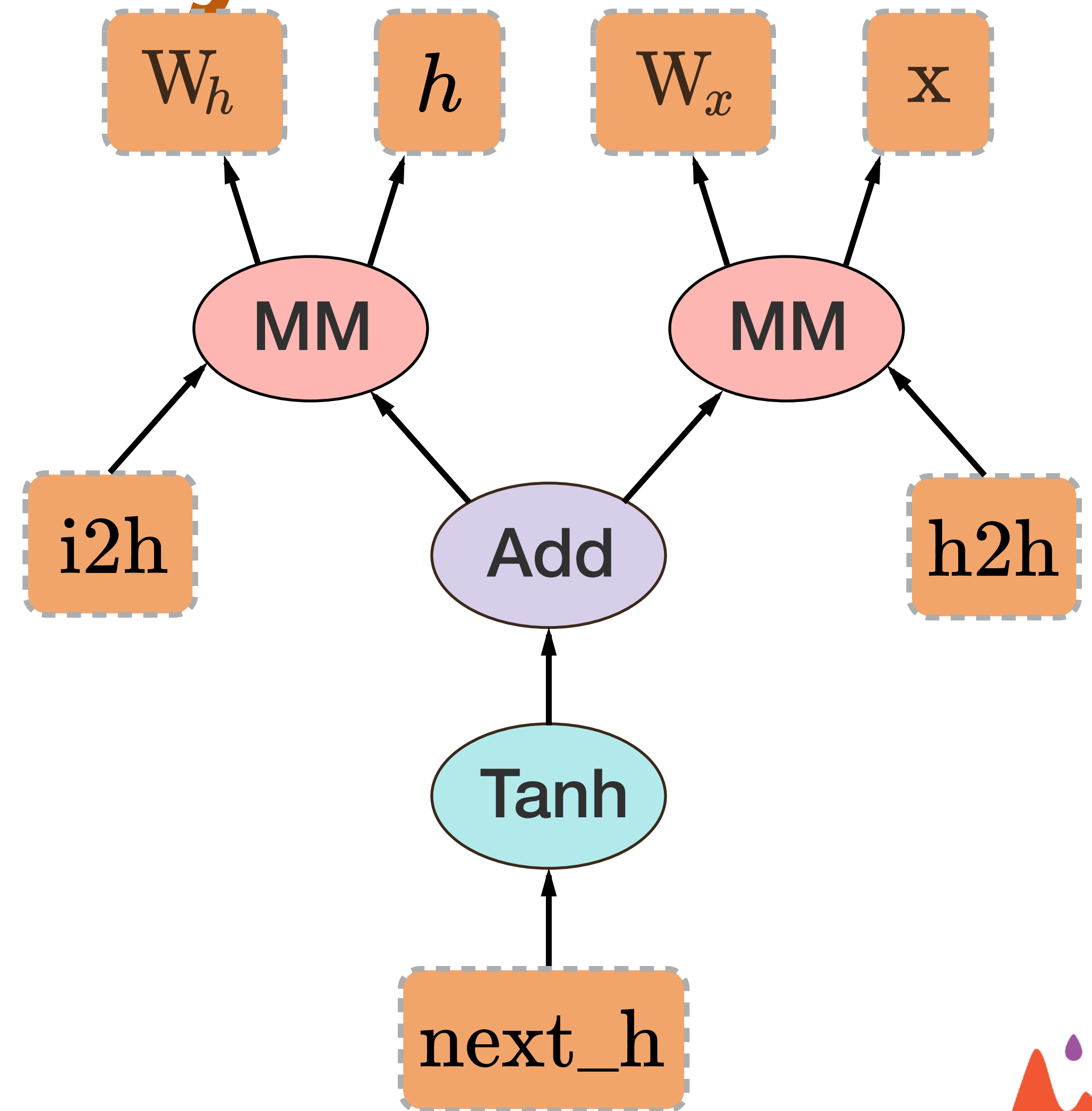


Imperative Frameworks: PyTorch

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```



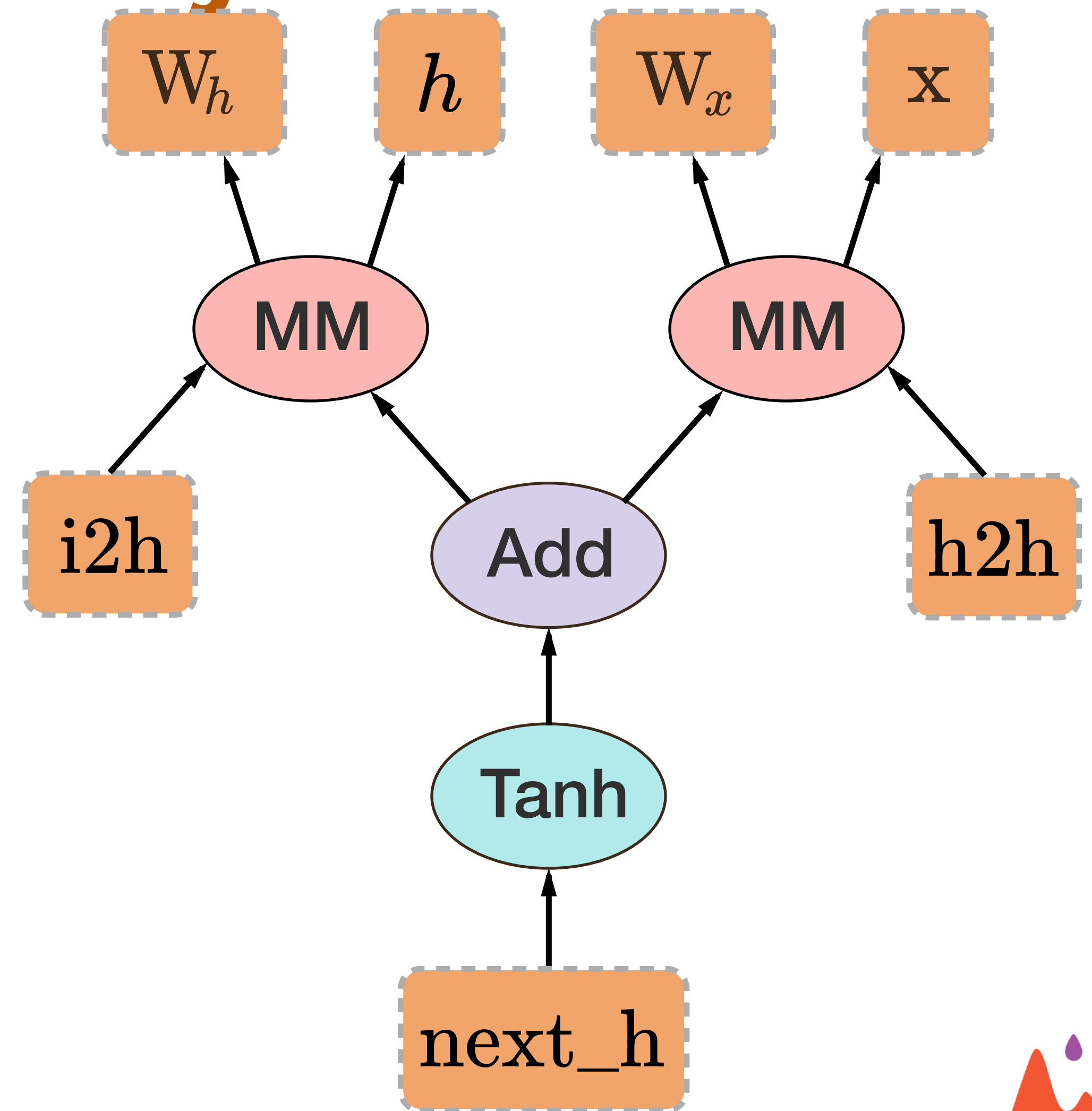
Imperative Frameworks: PyTorch

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```



Computation
Graphs

Deep Learning
Frameworks

Imperative
Frameworks

JIT
Compilation



Imperative Frameworks: PyTorch

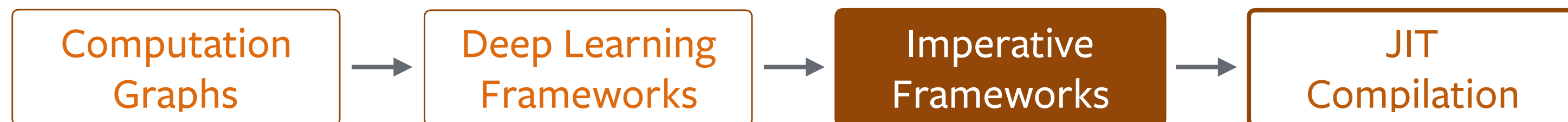
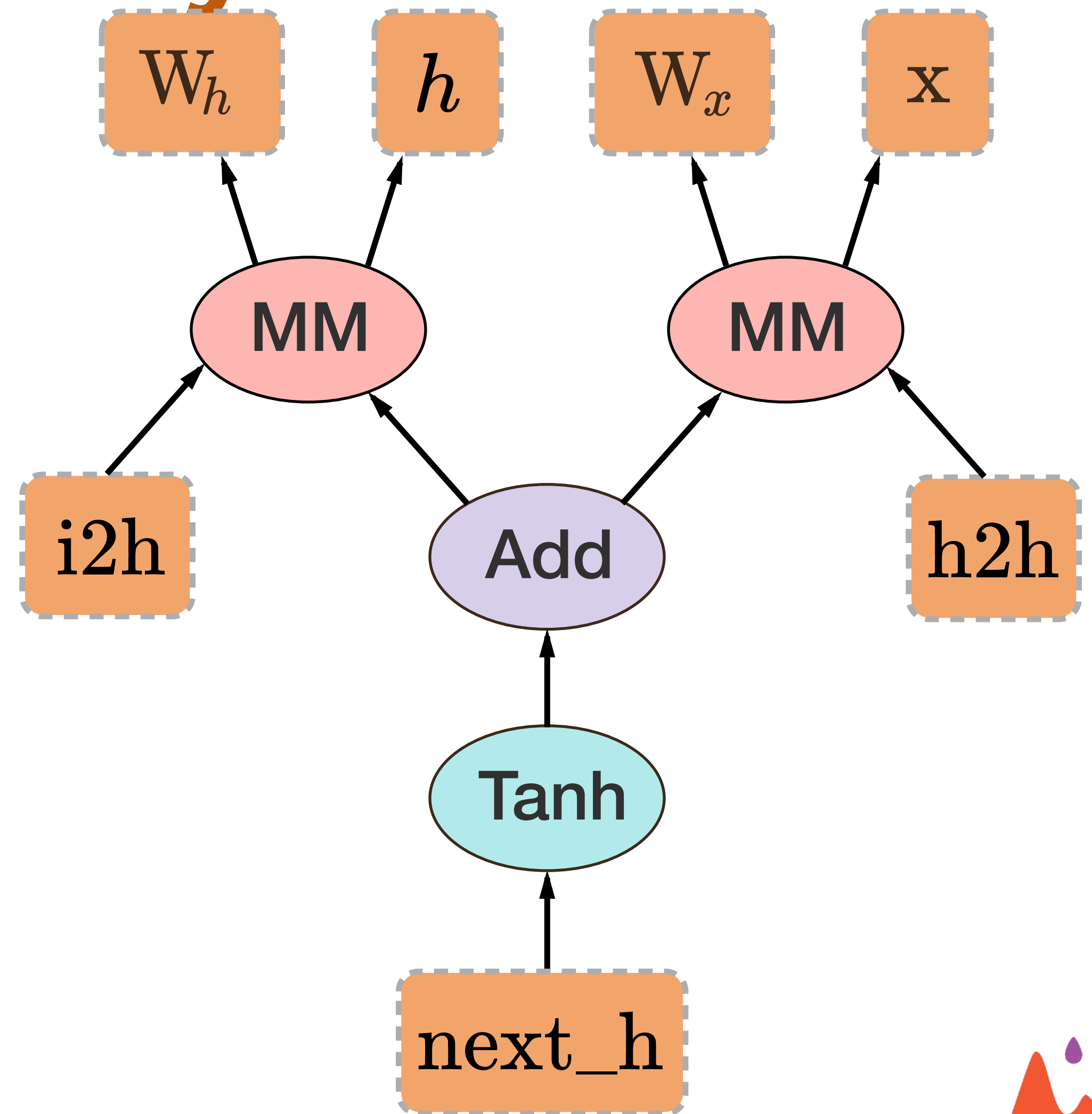
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```

Hence, graph construction
has to be FAST



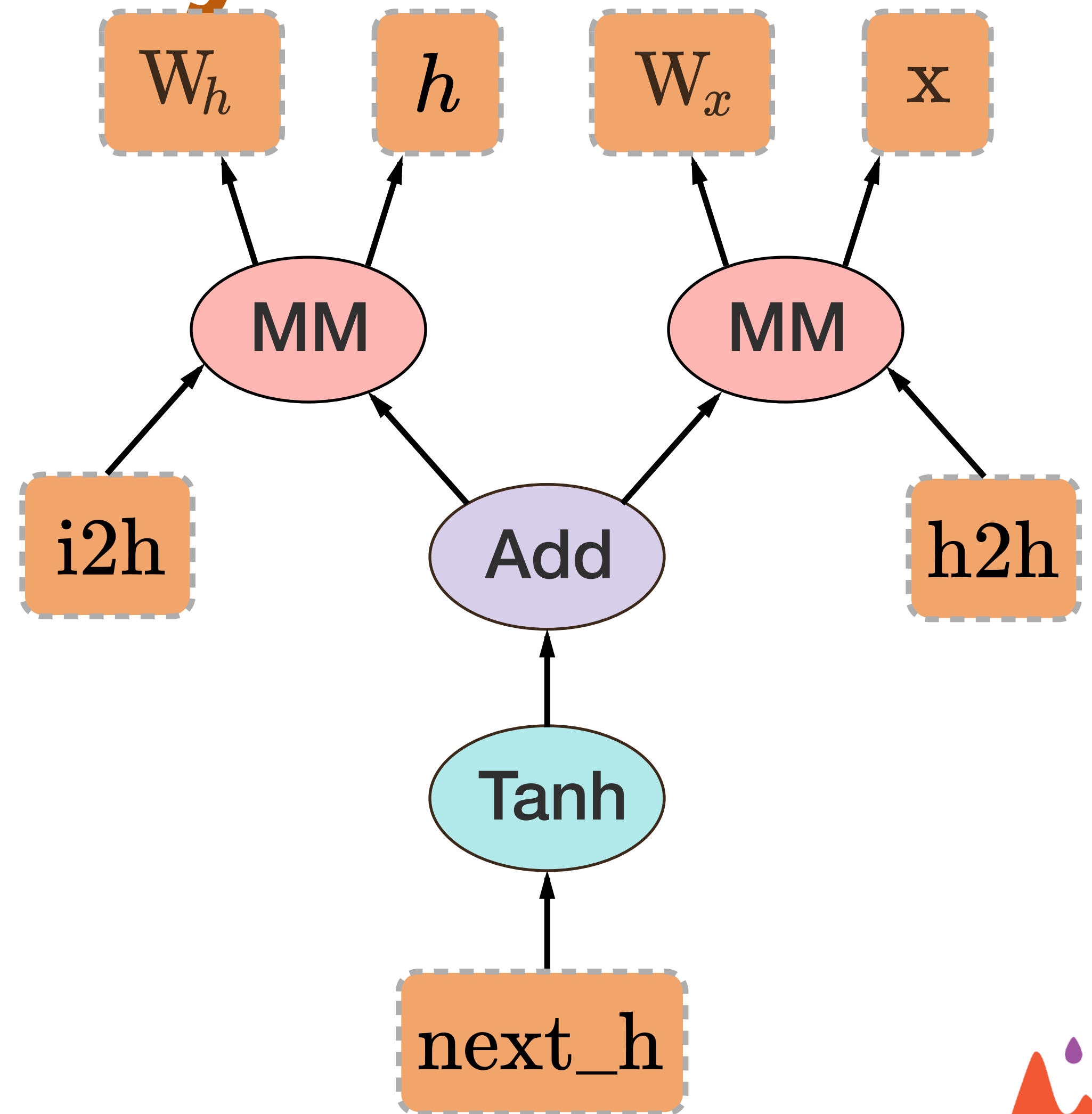
Imperative Frameworks: PyTorch

```
from torch.autograd import Variable
```

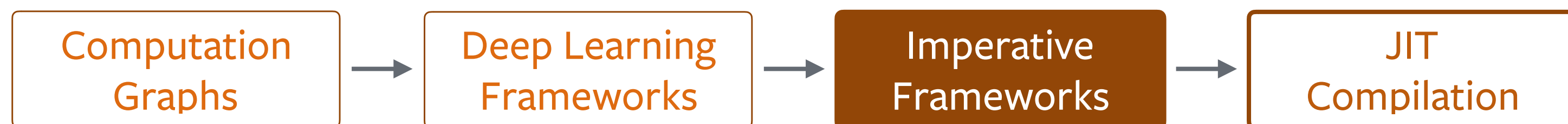
```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```



Initially written in Python



Imperative Frameworks: PyTorch

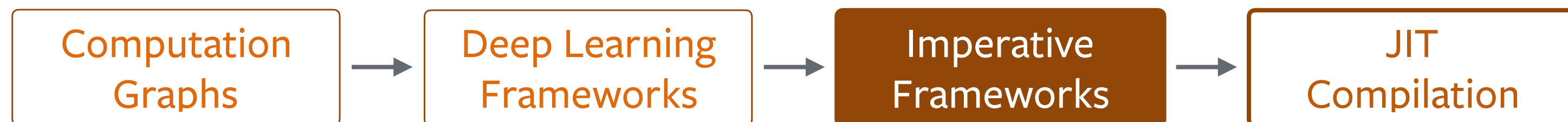
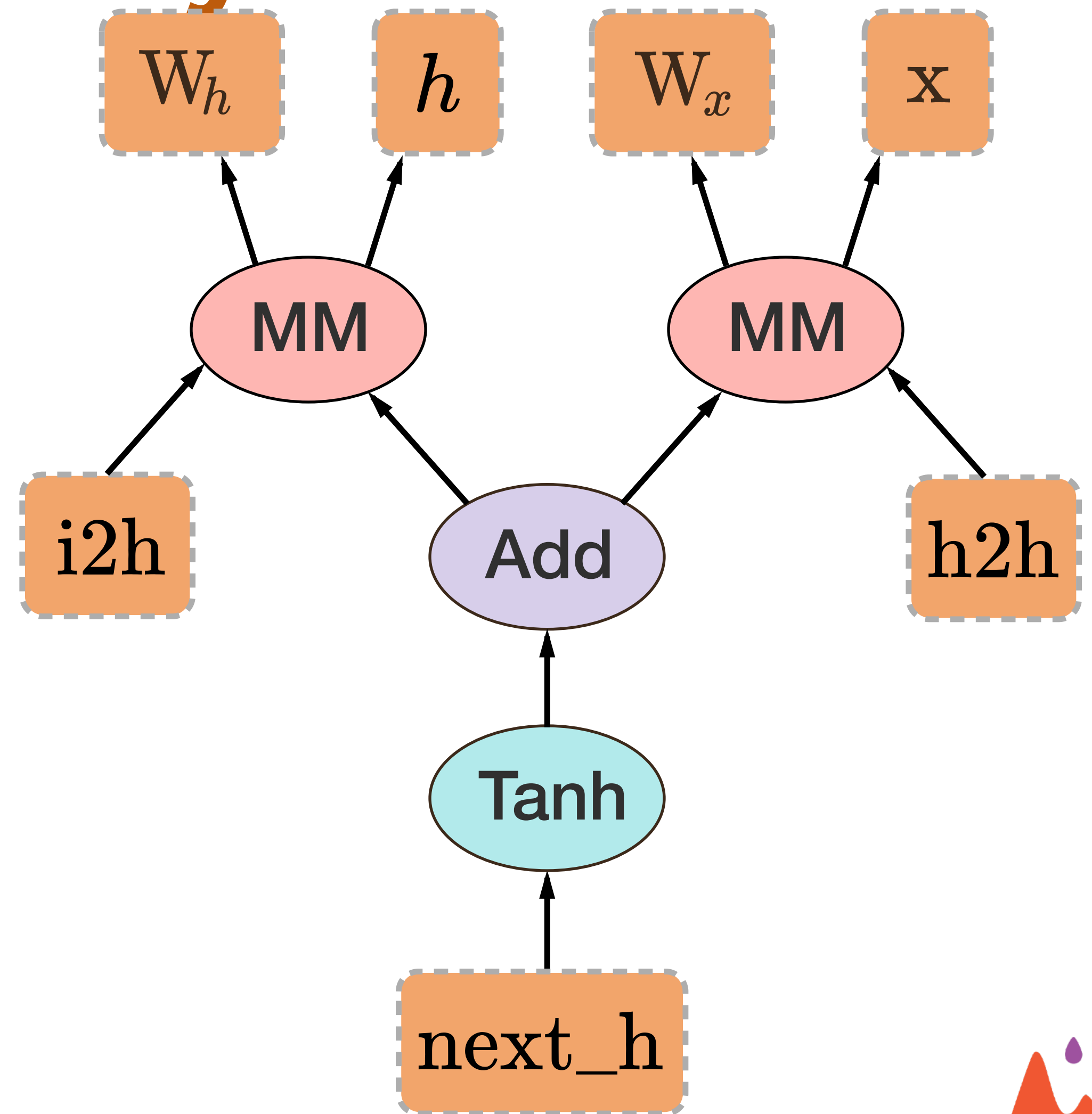
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```

Overhead too high



Imperative Frameworks: PyTorch

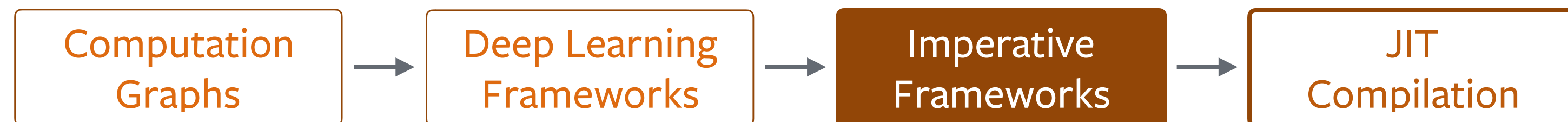
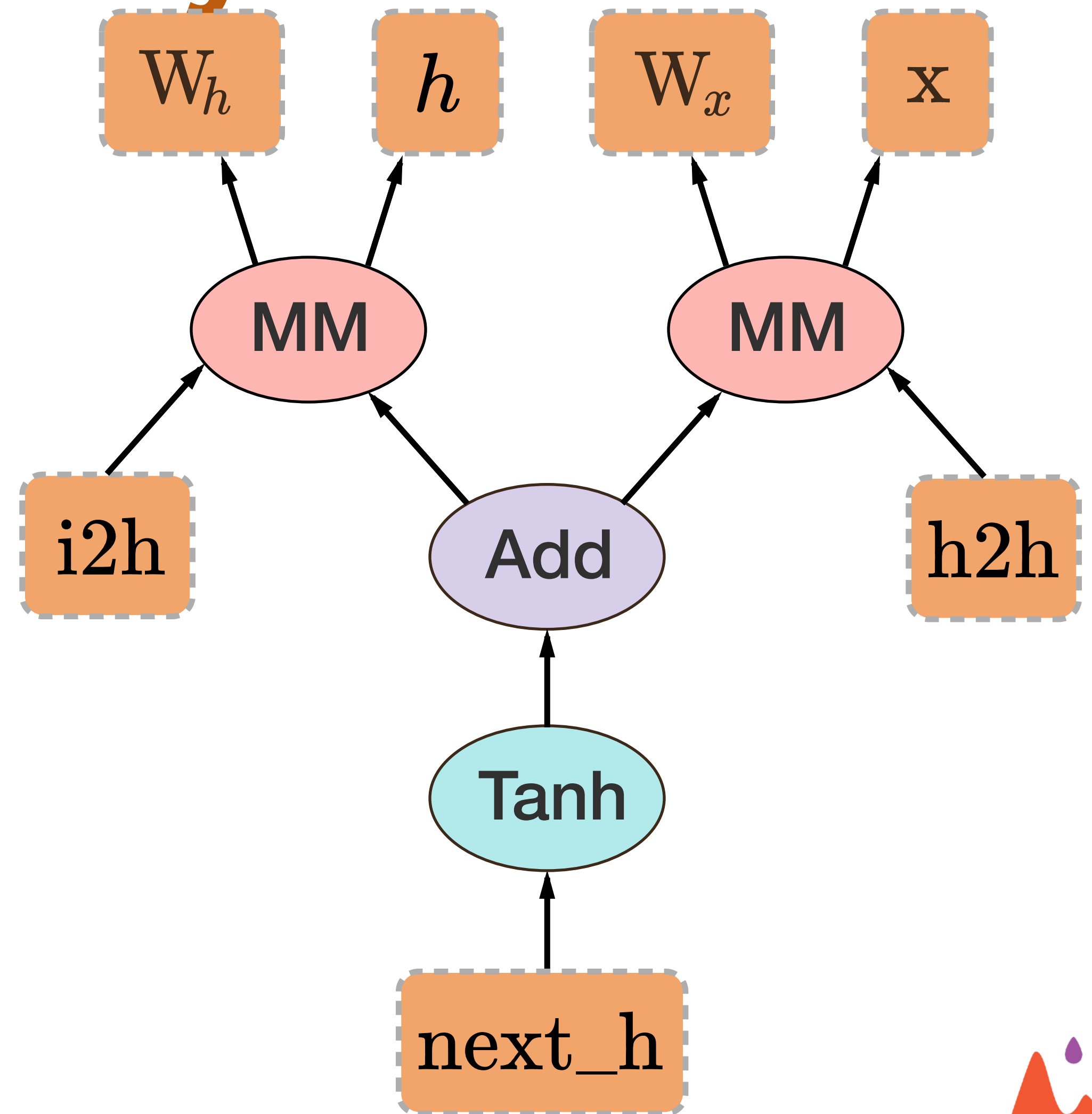
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```

Moved to CPython



Imperative Frameworks: PyTorch

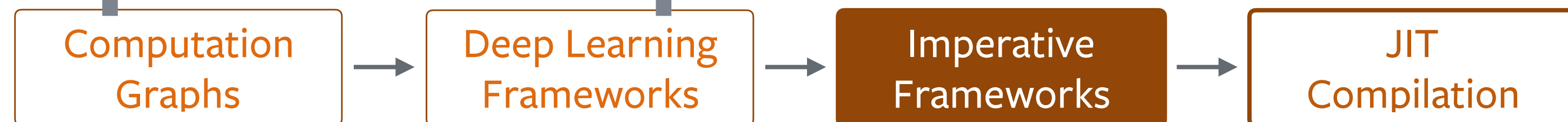
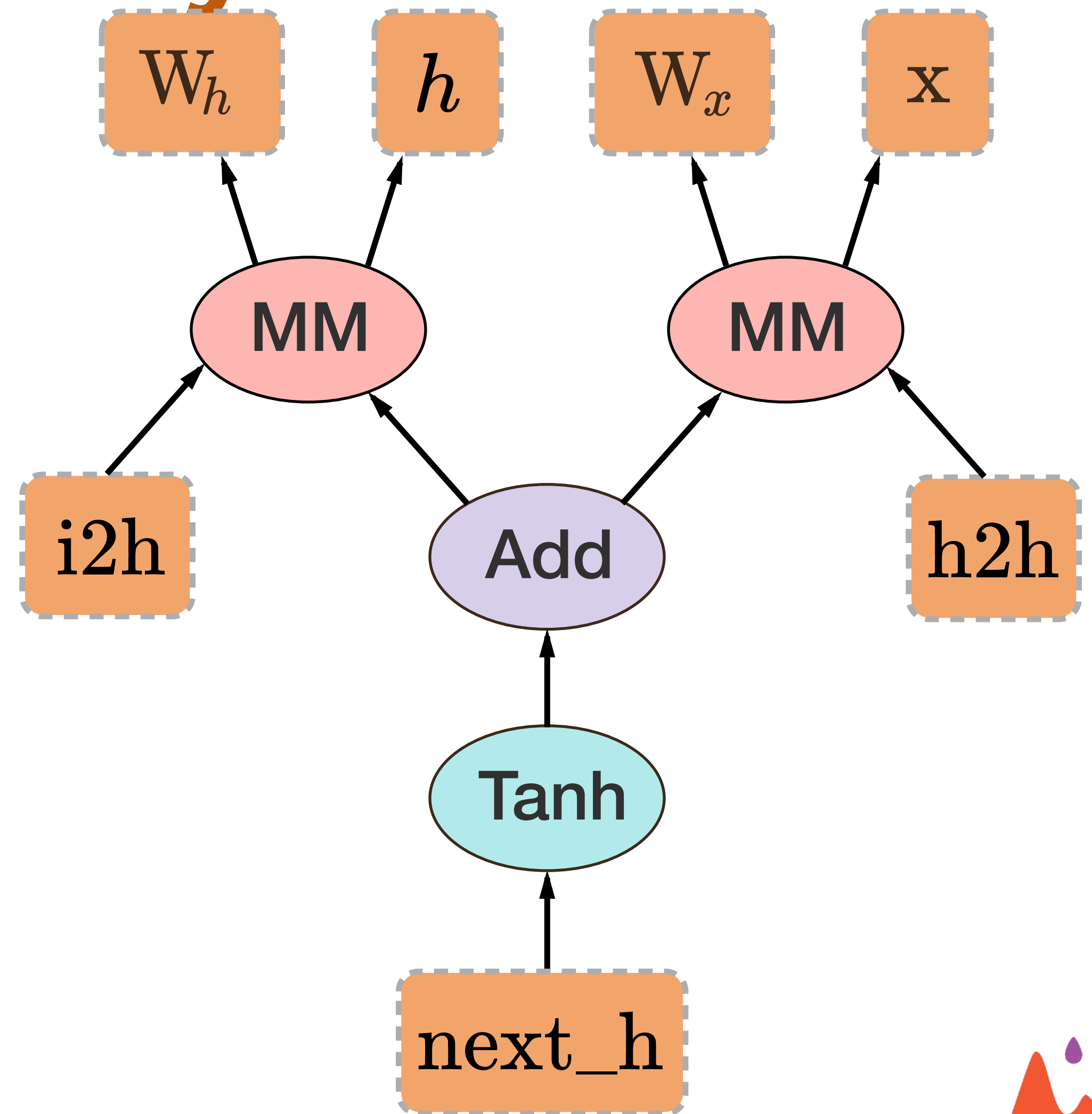
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```

Moved to CPython
used Flame graphs to find
and optimize hotspots



Imperative Frameworks: PyTorch

```
from torch.autograd import Variable
```

SNAKEVIZ

```
x = Variable(  
prev_h = Vari  
W_h = Variabl  
W_x = Variabl  
  
i2h = torch.n  
h2h = torch.n  
next_h = i2h  
next_h = next
```

SnakeViz

[Installation](#)

[Starting SnakeViz](#)

[Generating Profiles](#)

[Interpreting Results](#)

[Controls](#)

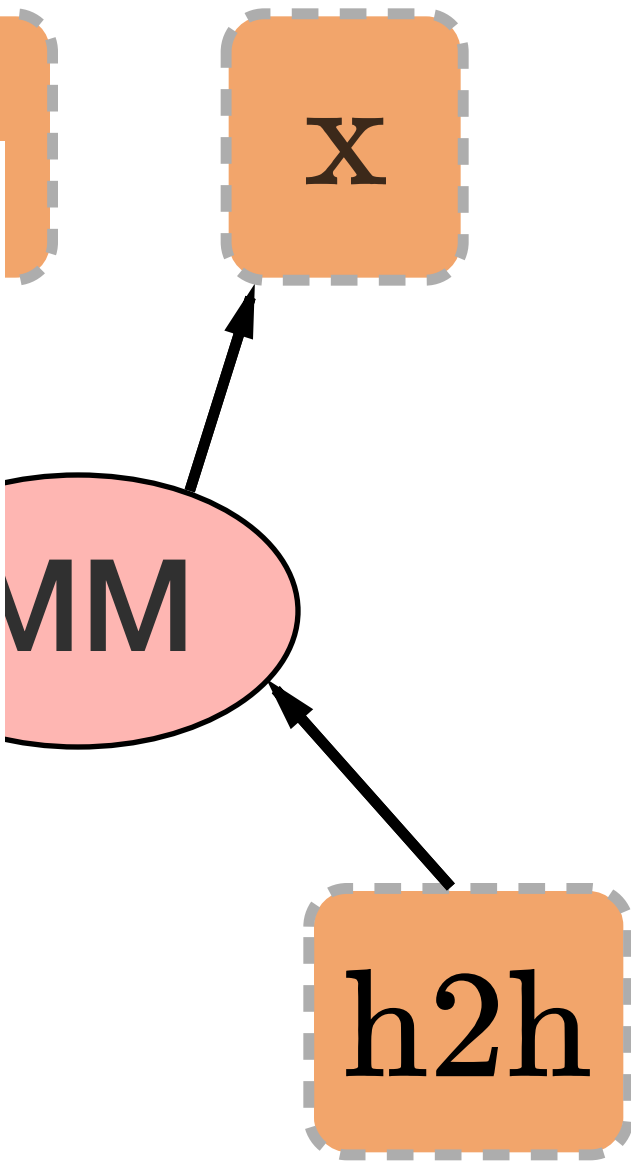
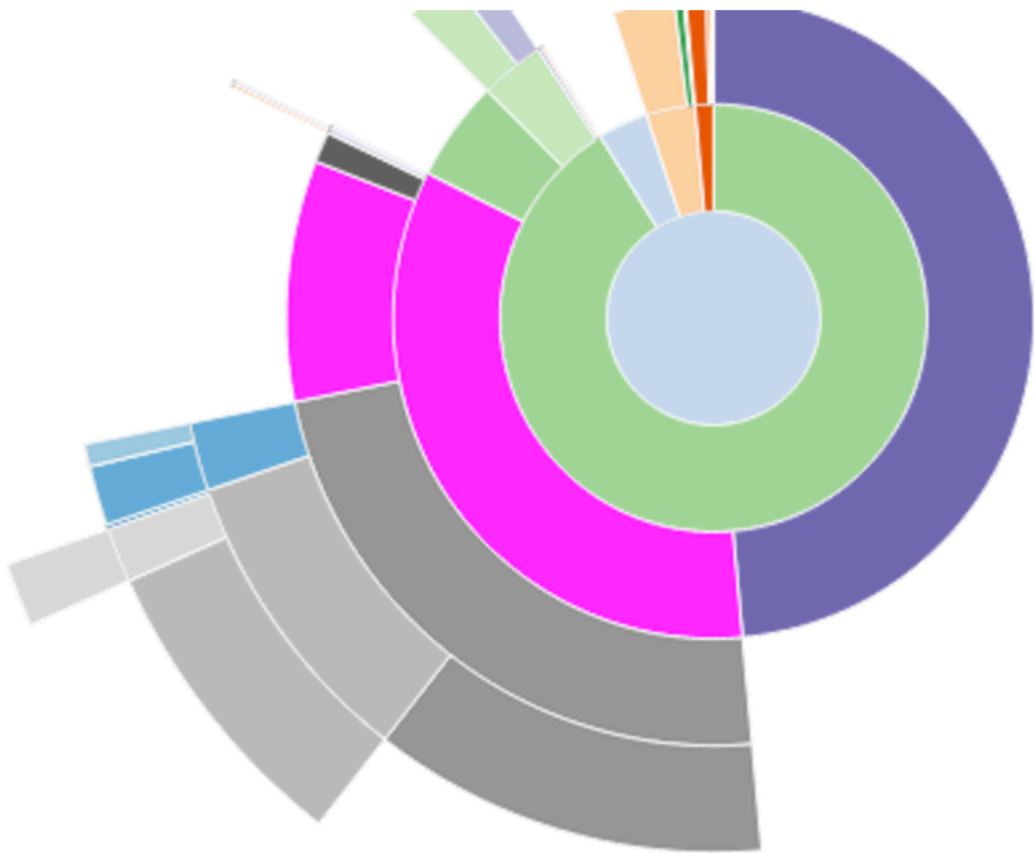
[Notes](#)

[Contact](#)

FUNCTION INFO

Placing your cursor over an arc will highlight that arc and any other visible instances of the same function call. It will also display a list of information to the left of the sunburst.

Name:
filter
Cumulative Time:
0.000294 s (31.78 %)
File:
fnmatch.py
Line:
48
Directory:
/Users/jiffyclub/miniconda3/en
vs/snakevizdev/lib/python3.4/

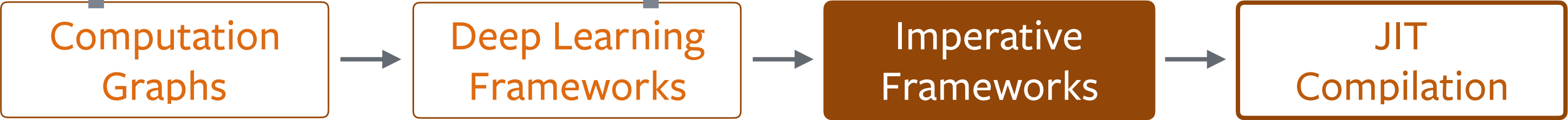
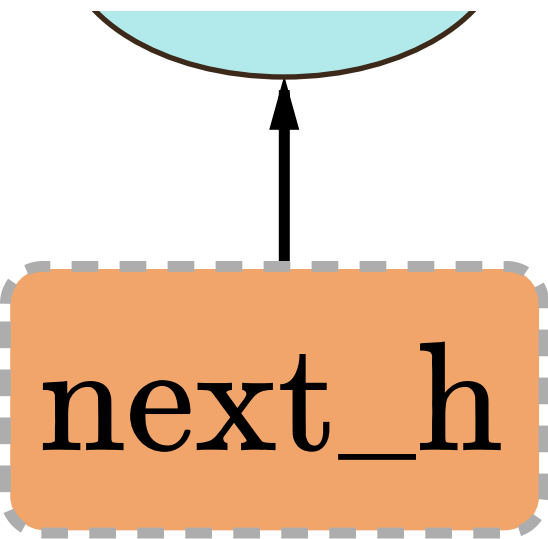


← PREVIOUS NEXT → [GITHUB](#)

```
next_h.backward(torch.ones(1, 20))
```

The displayed information includes:

Moved to CPython
used Flame graphs to find
and optimize hotspots



Imperative Frameworks: PyTorch

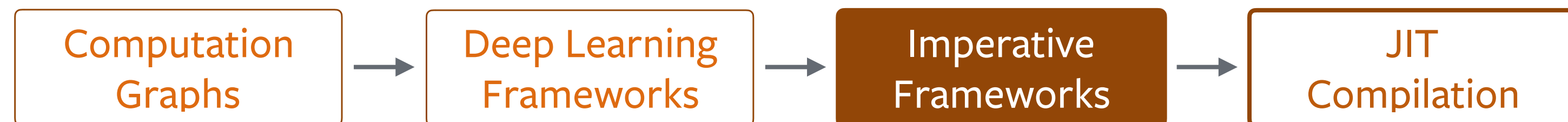
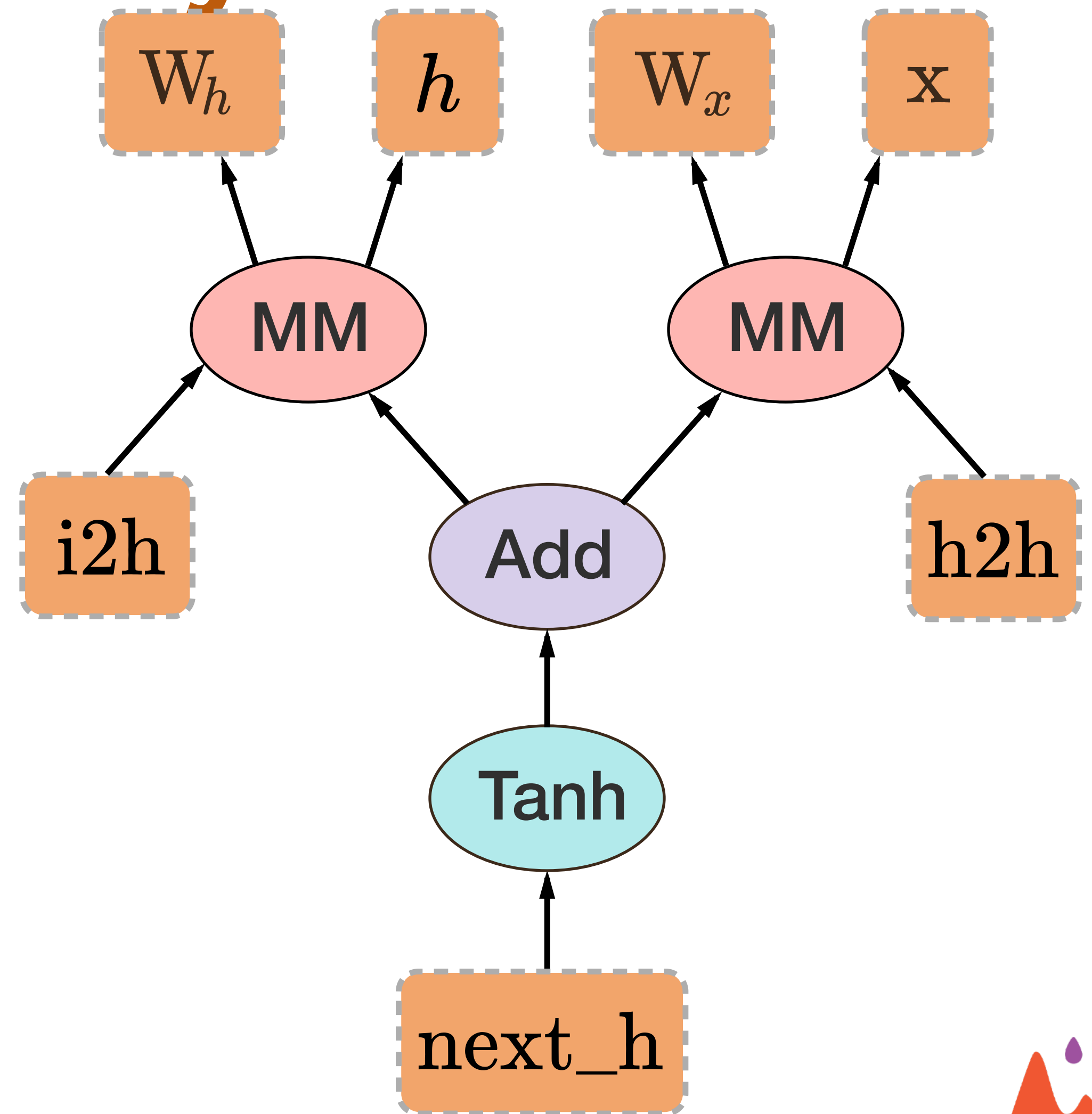
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```

**Overall speed as fast as
declarative frameworks**



Performance stats from Nov 2016

Task	Torch		PyTorch	
ResNet-101	544ms	10GB (5,4GB)	516ms	4,9GB
ResNet-101 2GPU	580ms	10GB (5,6GB)	640ms	4,9GB
Penn Treebank 2-layer LSTM	57ms	865MB	62ms	370MB



Performance stats from Nov 2016

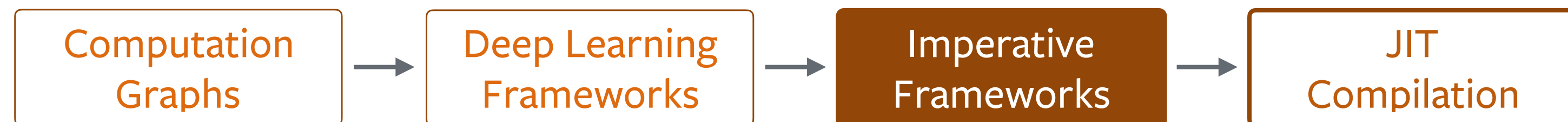
Task	Torch		PyTorch	
ResNet-101	544ms	10GB (5,4GB)	516ms	4,9GB
ResNet-101 2GPU	580ms	10GB (5,6GB)	640ms	4,9GB
Penn Treebank 2-layer LSTM	57ms	865MB	62ms	370MB

Pure Python internals = **140ms**



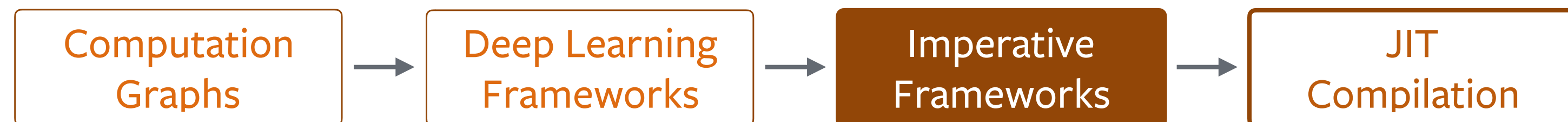
Graph Construction Speed

- PyTorch: nanoseconds to microseconds
- TensorFlow: milliseconds to several seconds
- Theano: seconds to minutes (hours?)
- MXNet: i dont know -> ask your instructor :)



Additional valuable features in PyTorch

- Low memory usage even without a static optimizer



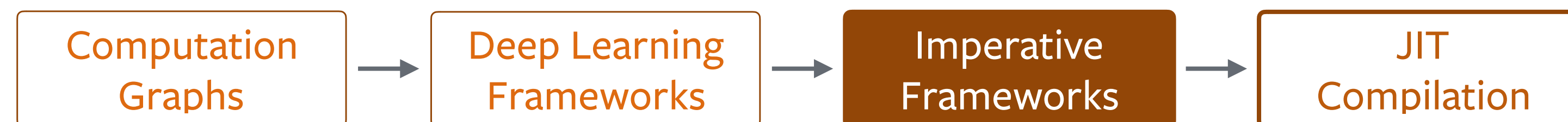
Performance stats from Nov 2016

Task	Torch		PyTorch	
ResNet-101	544ms	10GB (5,4GB)	516ms	4,9GB
ResNet-101 2GPU	580ms	10GB (5,6GB)	640ms	4,9GB
Penn Treebank 2-layer LSTM	57ms	865MB	62ms	370MB



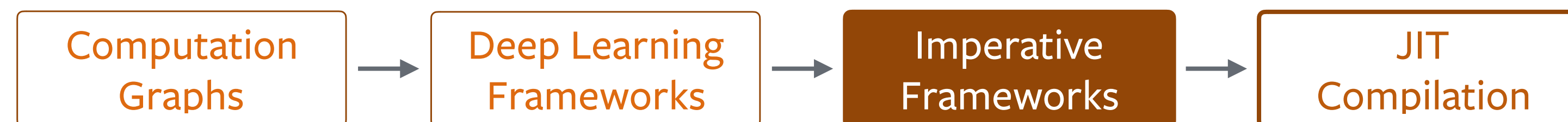
Additional valuable features in PyTorch

- Low memory usage even without a static optimizer
- Intermediate buffers are always freed
- Developers given constructs to allocate temporary buffers
 - save_for_backward
 - requires_grad



Additional valuable features in PyTorch

- Low memory usage even without a static optimizer
- Intermediate buffers are always freed
- Developers given constructs to allocate temporary buffers
 - save_for_backward
 - requires_grad
- in-place operations

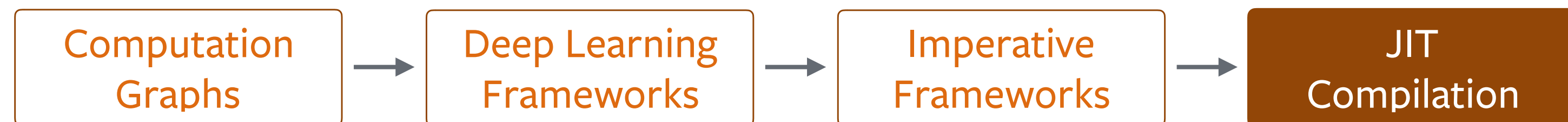


JIT Compilation



JIT Compilation

- Possible in Imperative Frameworks
- The key idea is deferred or lazy evaluation
 - $y = x + 2$
 - $z = y * y$
 - # nothing is executed yet, but the graph is being constructed
 - `print(z)` # now the entire graph is executed: $z = (x+2) * (x+2)$
- We can do just in time compilation on the graph before execution



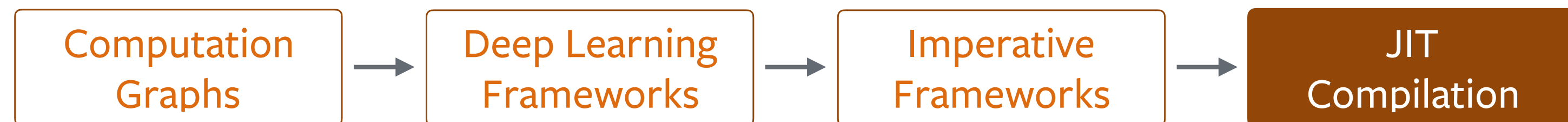
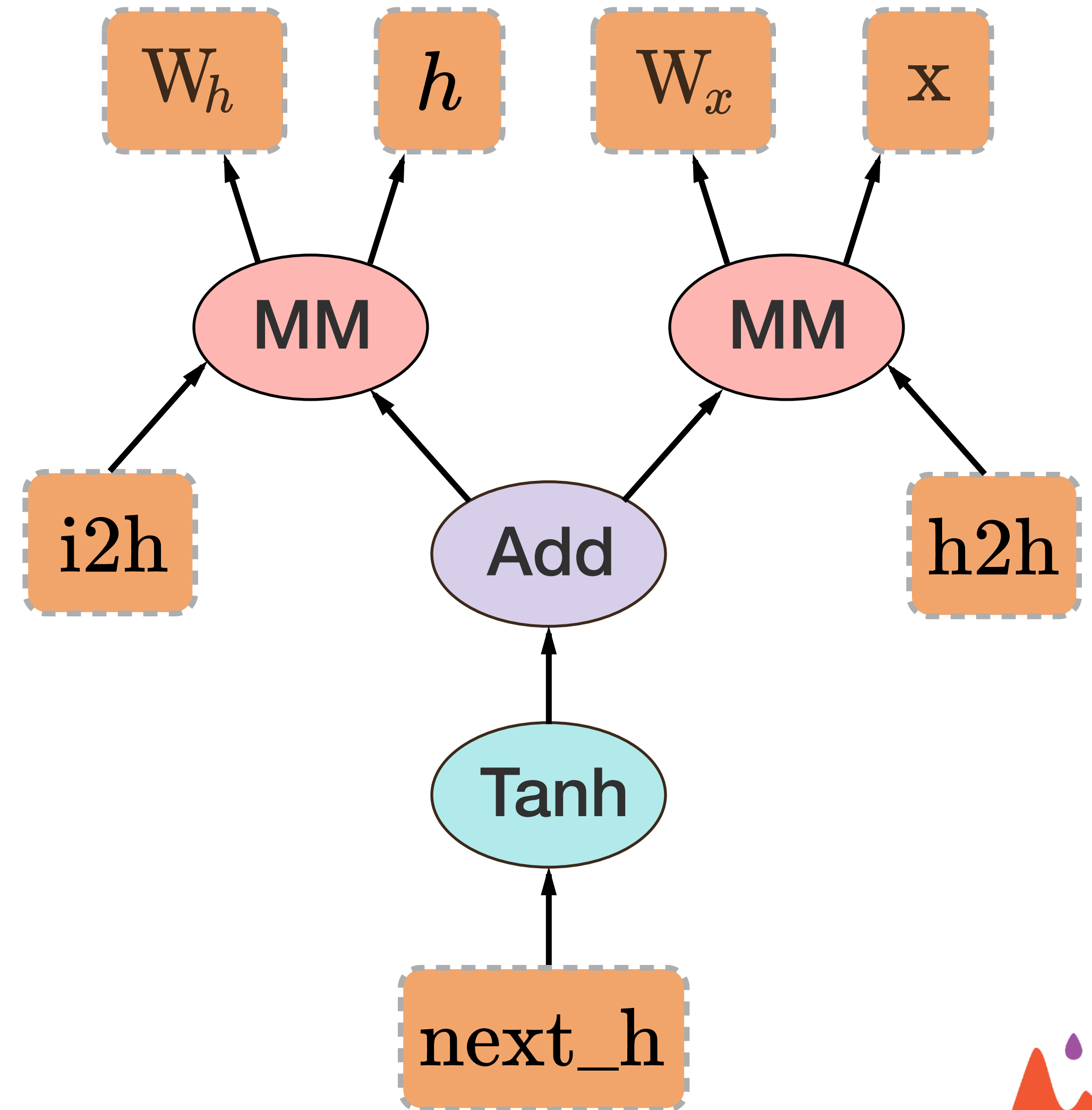
Lazy Evaluation

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```



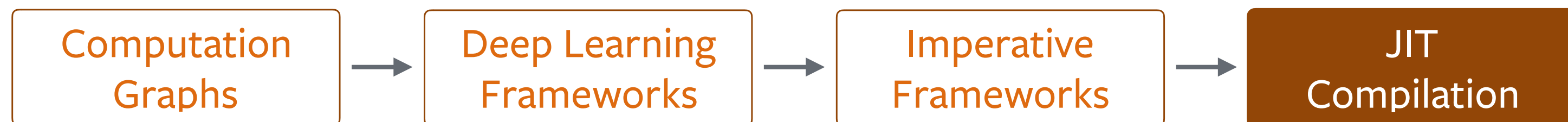
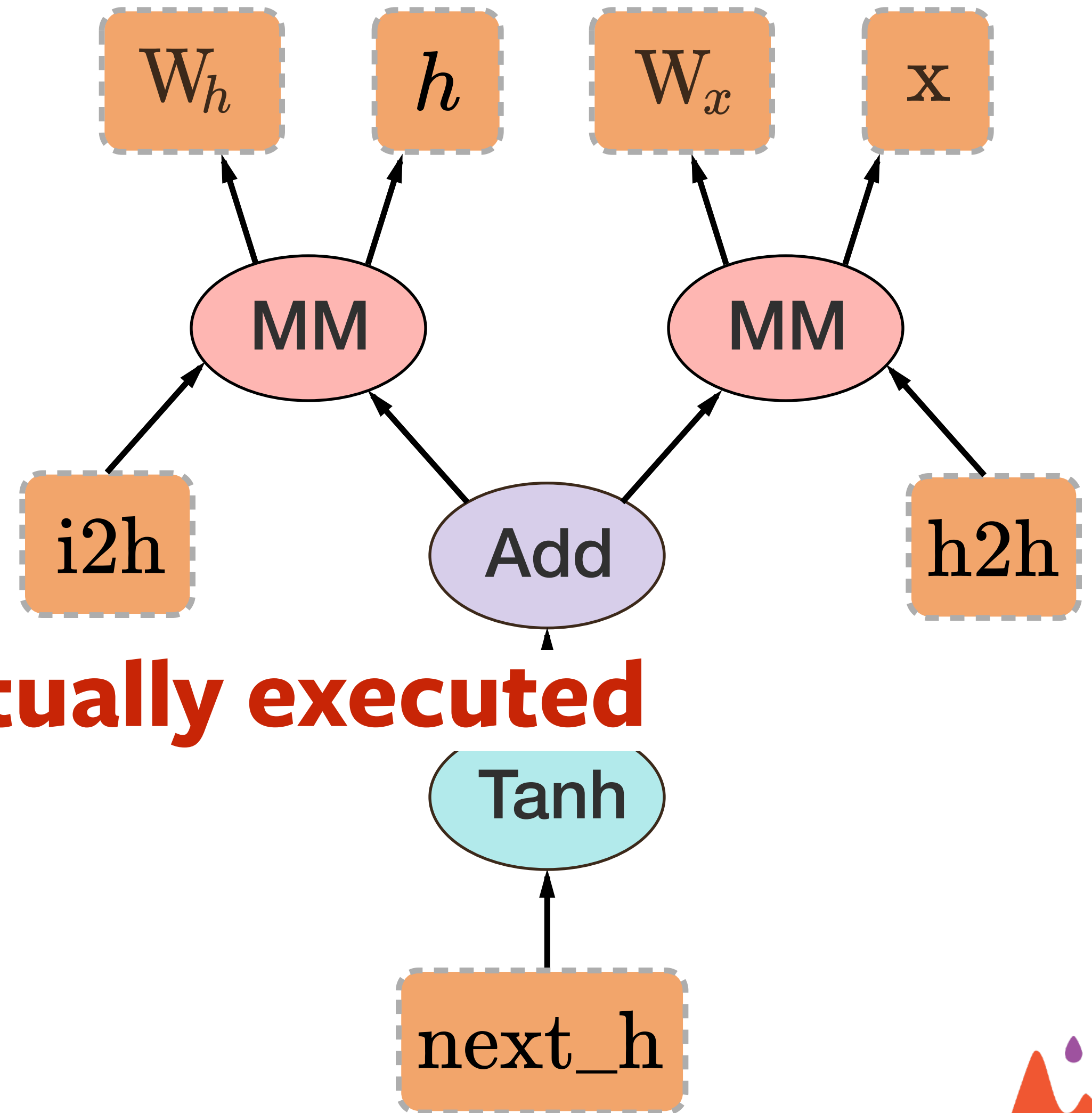
Lazy Evaluation

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

Graph built but not actually executed



Lazy Evaluation

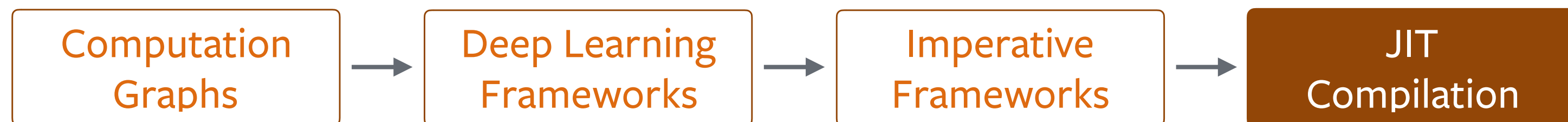
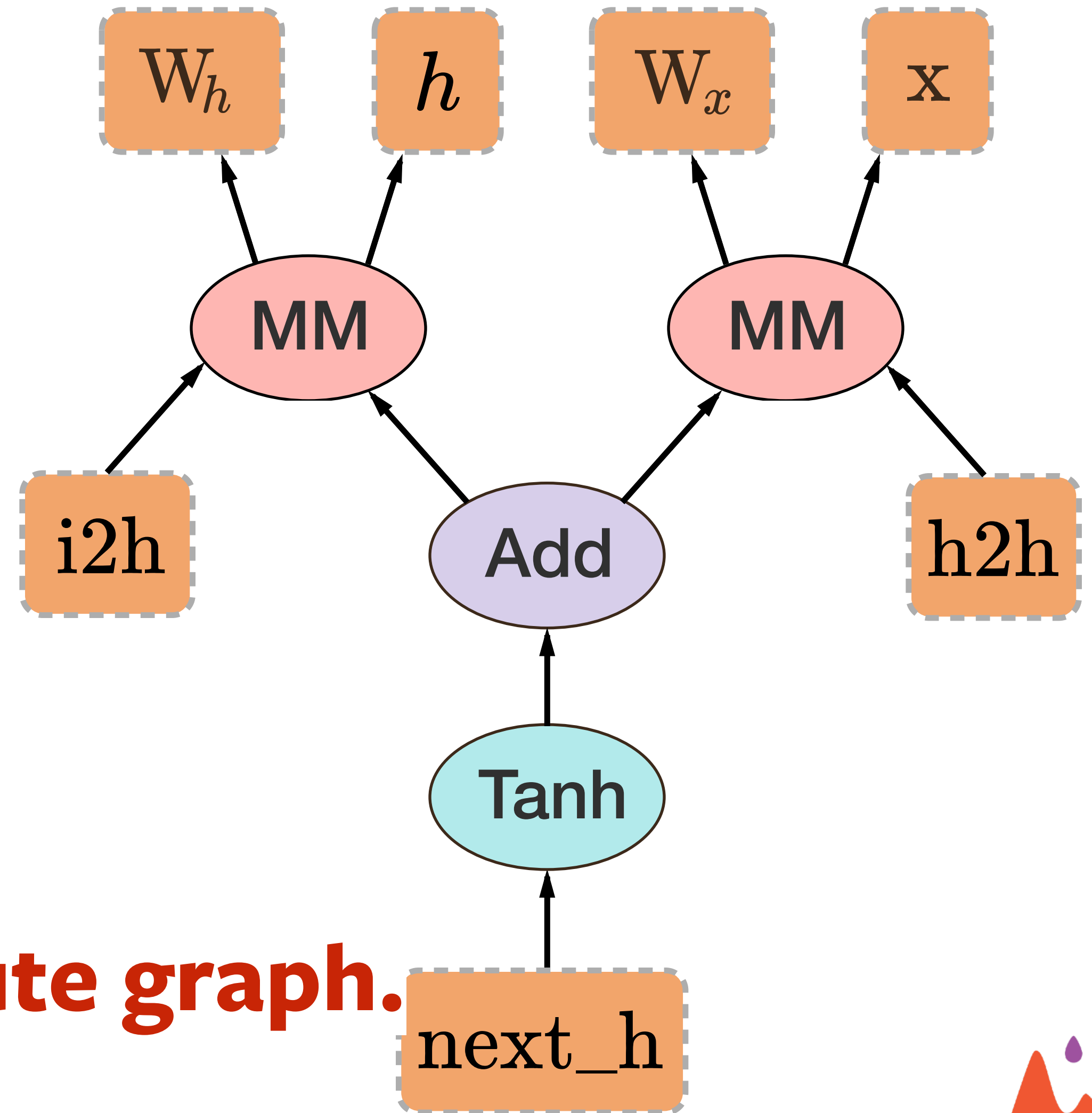
```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

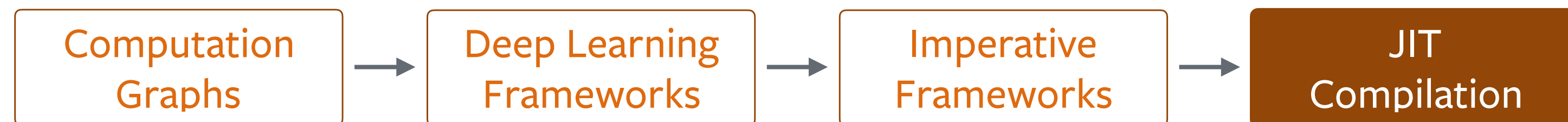
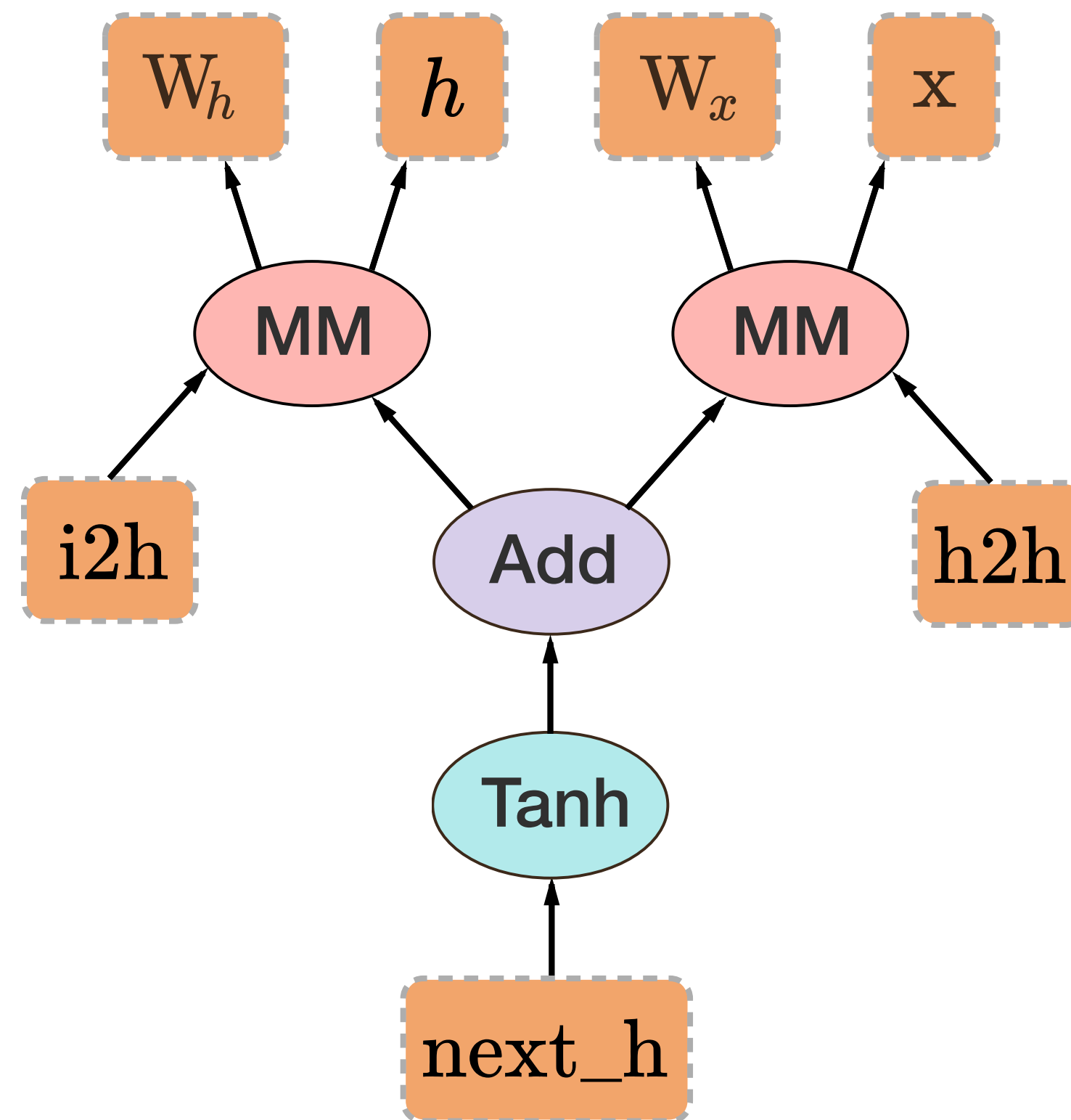
```
print(next_h)
```

Data accessed. Execute graph.



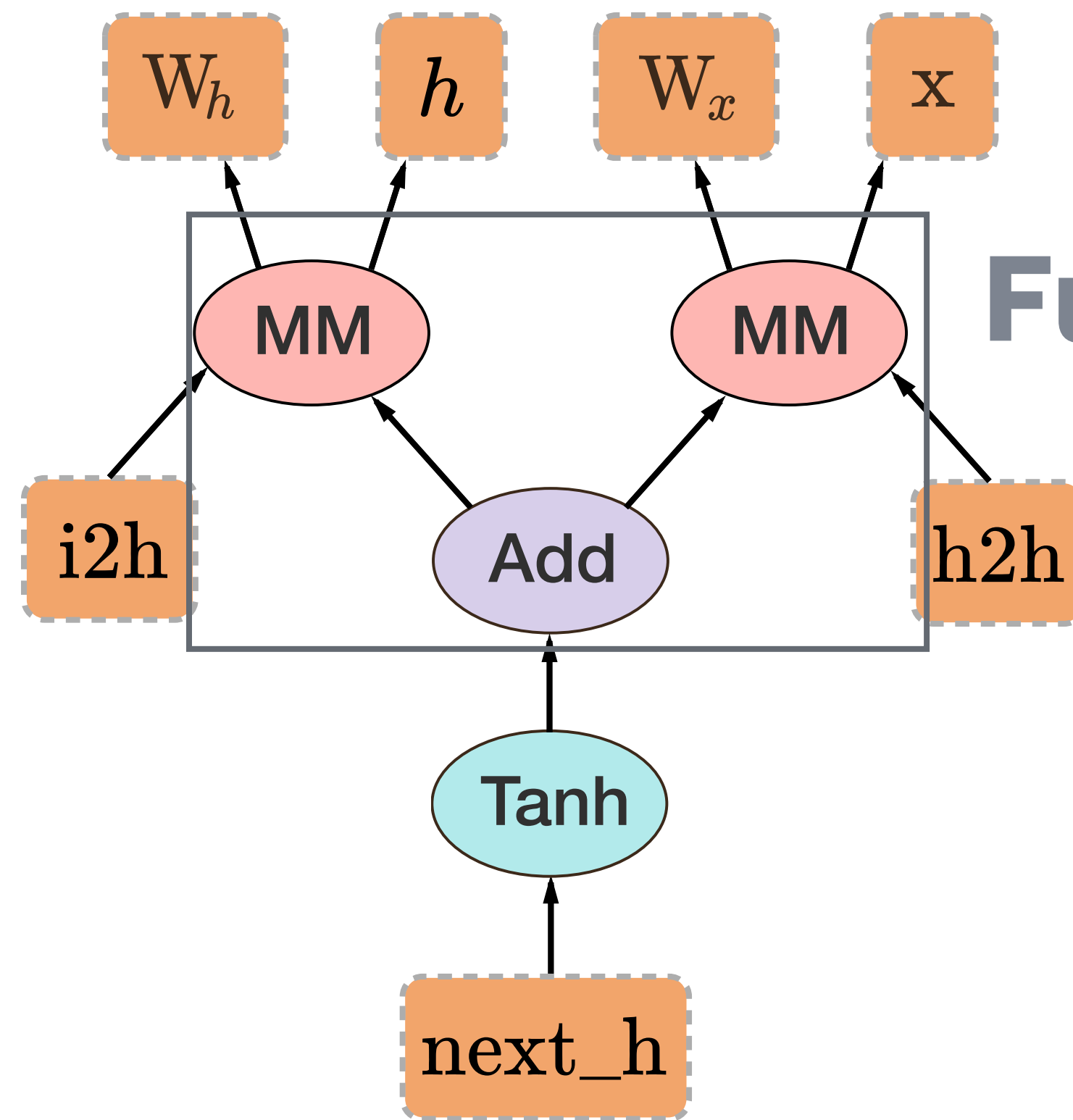
Lazy Evaluation

- A little bit of time between building and executing graph
 - Use it to compile the graph just-in-time



JIT Compilation

- Fuse and optimize operations



Fuse operations. Ex:

```
1  x = [0, 1, 2, 3, 4]
2  for i in range(len(x)):
3      x[i] = x[i] + 1
4
5  for i in range(len(x)):
6      x[i] = x[i] * 2
7
8
9  # Fused
10 for i in range(len(x)):
11     x[i] = (x[i] + 1) * 2
```



Computation
Graphs

Deep Learning
Frameworks

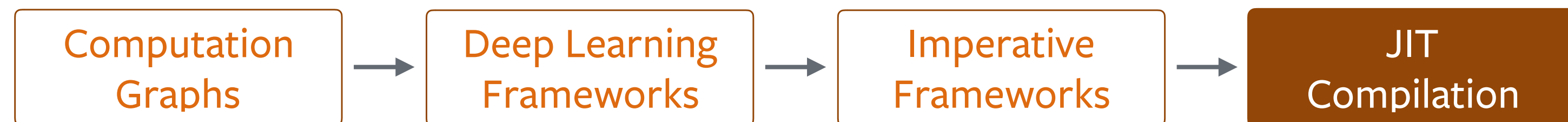
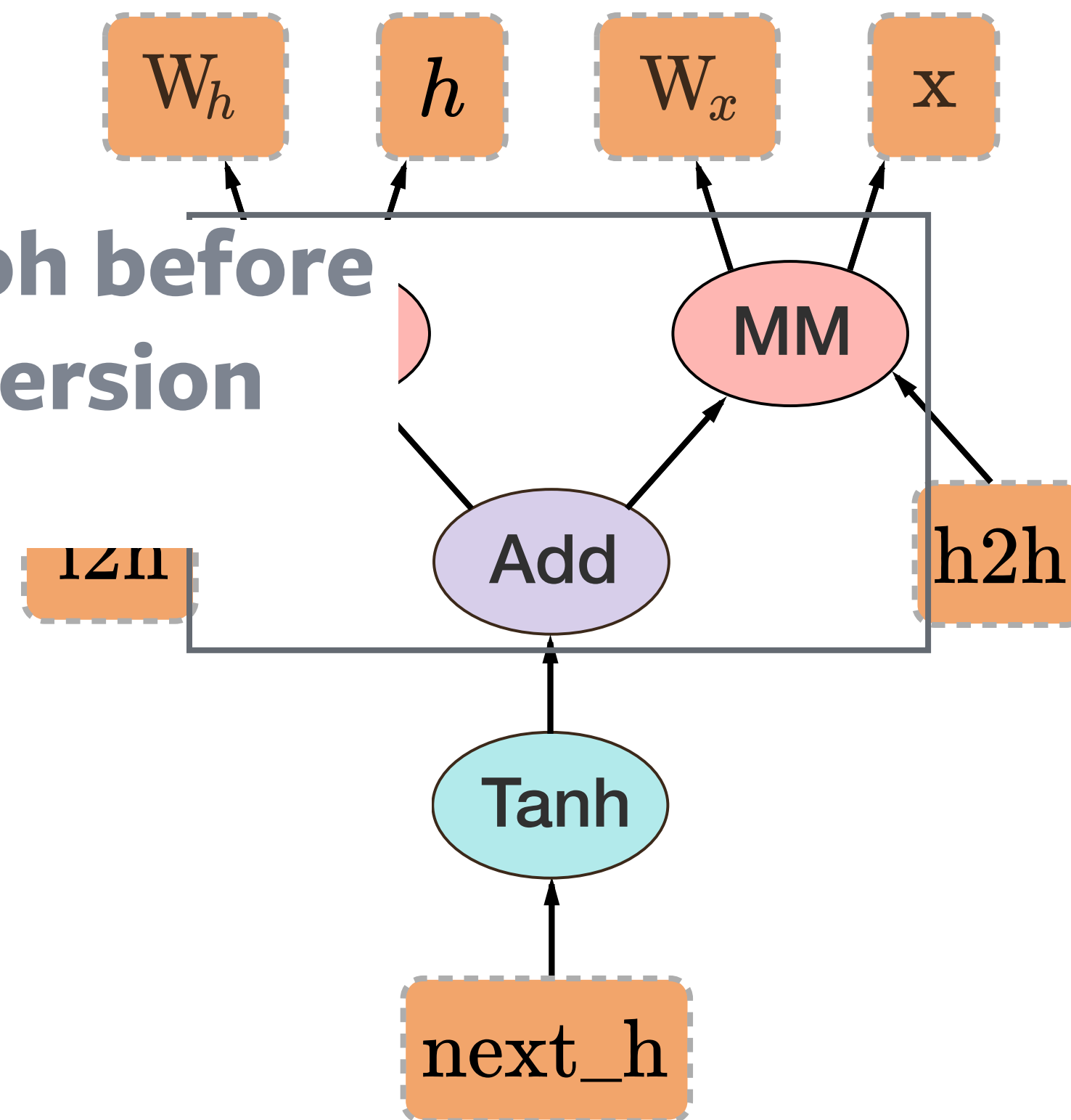
Imperative
Frameworks

JIT
Compilation

JIT Compilation

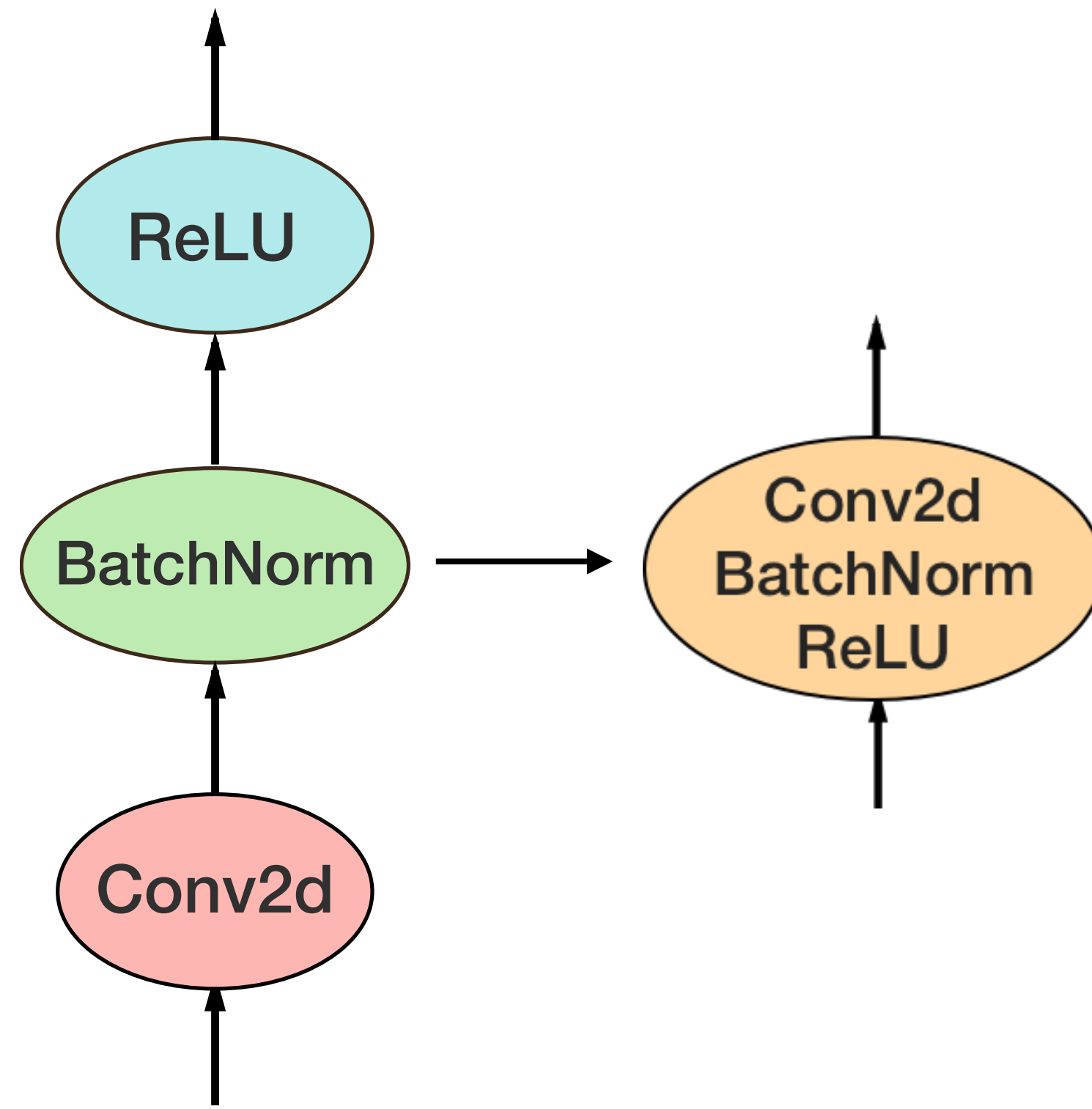
- Cache subgraphs

I've seen this part of the graph before
let me pull up the compiled version
from cache

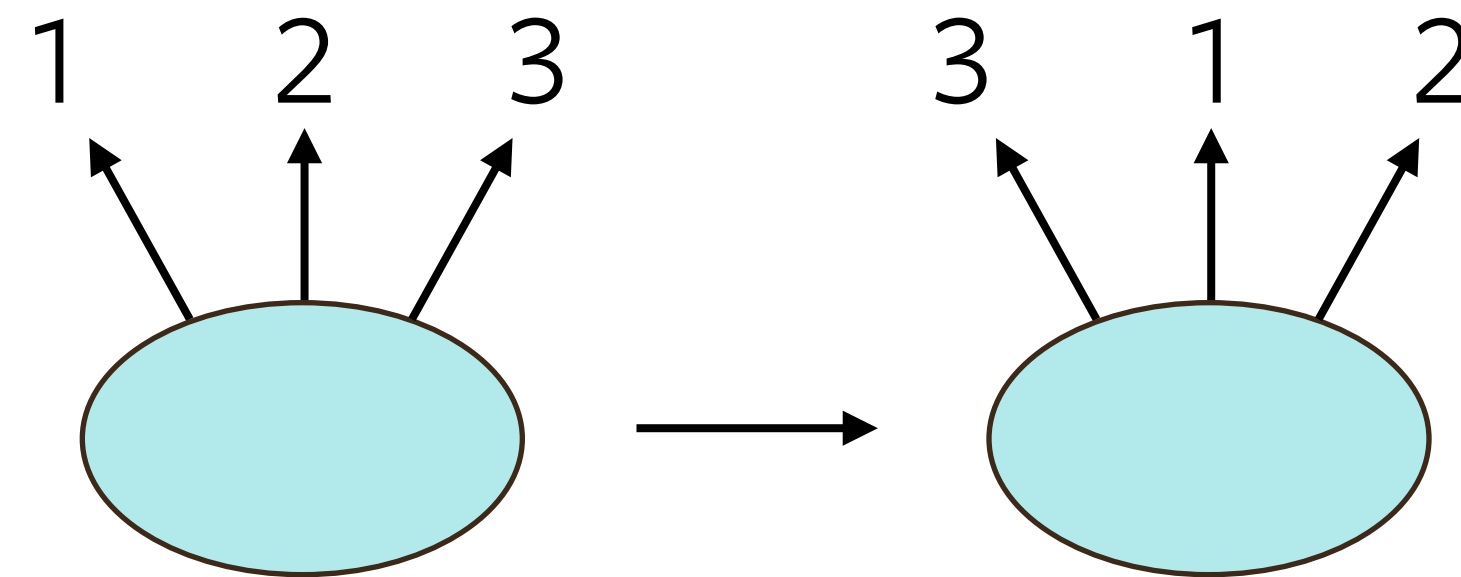


Compilation benefits

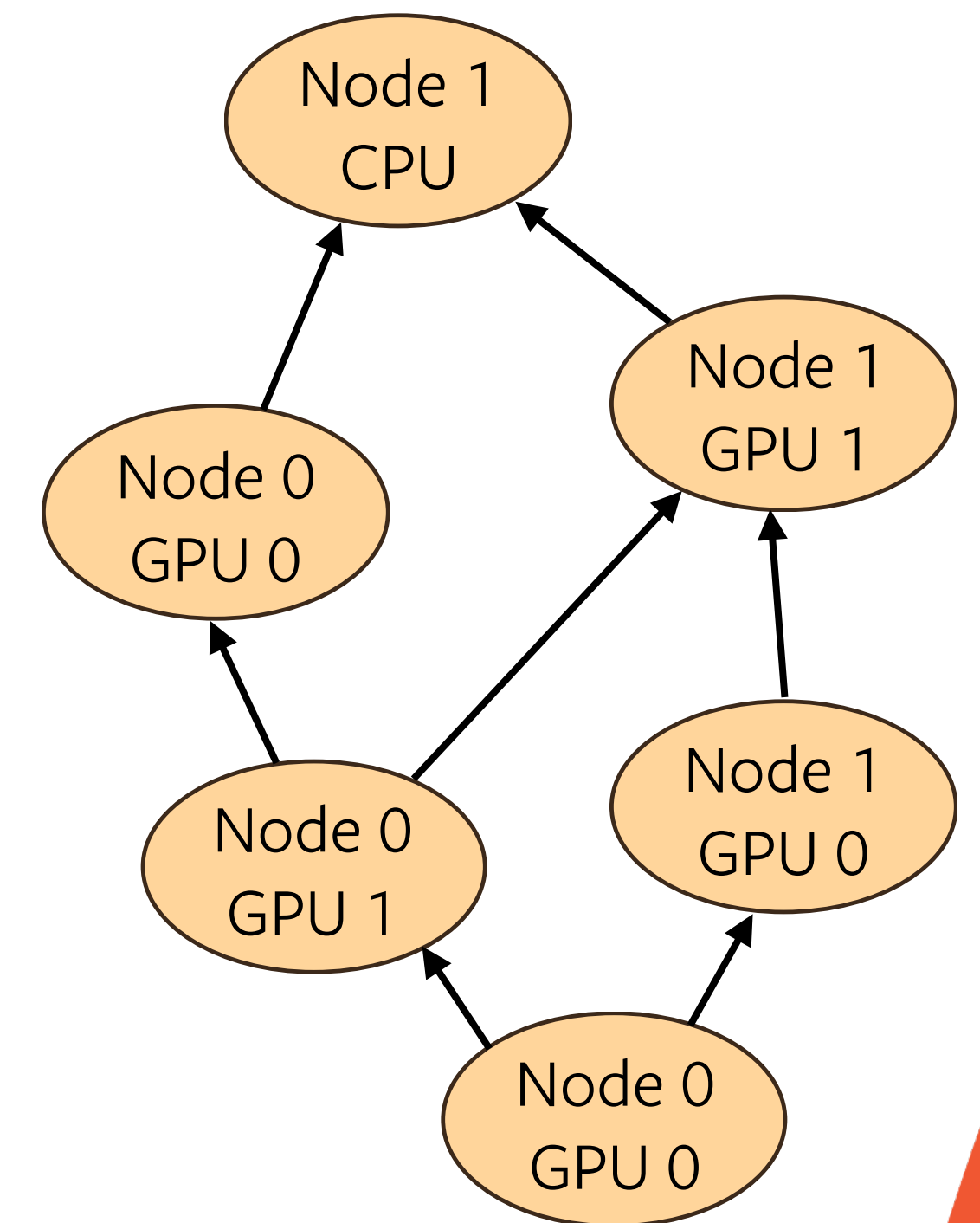
Kernel fusion



Out-of-order execution

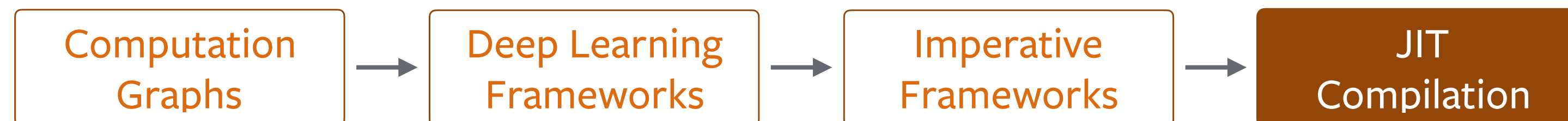


Automatic work placement



JIT Compilation

- Possible in Dynamic Frameworks
- The key idea is deferred or lazy evaluation
 - $y = x + 2$
 - $z = y * y$
 - # nothing is executed yet, but the graph is being constructed
 - `print(z)` # now the entire graph is executed: $z = (x+2) * (x+2)$
- We can do just in time compilation on the graph before execution
- We can cache repeating patterns in subsets of the graph
 - to avoid recompilation
- Compiler is very different from Ahead-of-time compiler
 - fast compilation
 - compile traces rather than full graph



Review

Computation
Graph Toolkits

Deep Learning
Frameworks

Imperative
Frameworks

JIT Compilation

Declarative

Imperative

Added features
on top of
computation graphs

Design
Challenges
Overheads

Lazy Execution
JIT Fusion
Subgraph caching
differences from AOT

Implementation

Advantages & Disadvantages



PYTORCH

With ❤️ from

<http://pytorch.org>

Released Jan 18th

40,000+ downloads

250+ community repos

4200+ user posts

330k+ forum views

facebook



You?