



**人工智能系统 System for AI**

**深度神经网络计算框架基础**

**Computation frameworks for DNN**

# 主要内容

- 主要内容

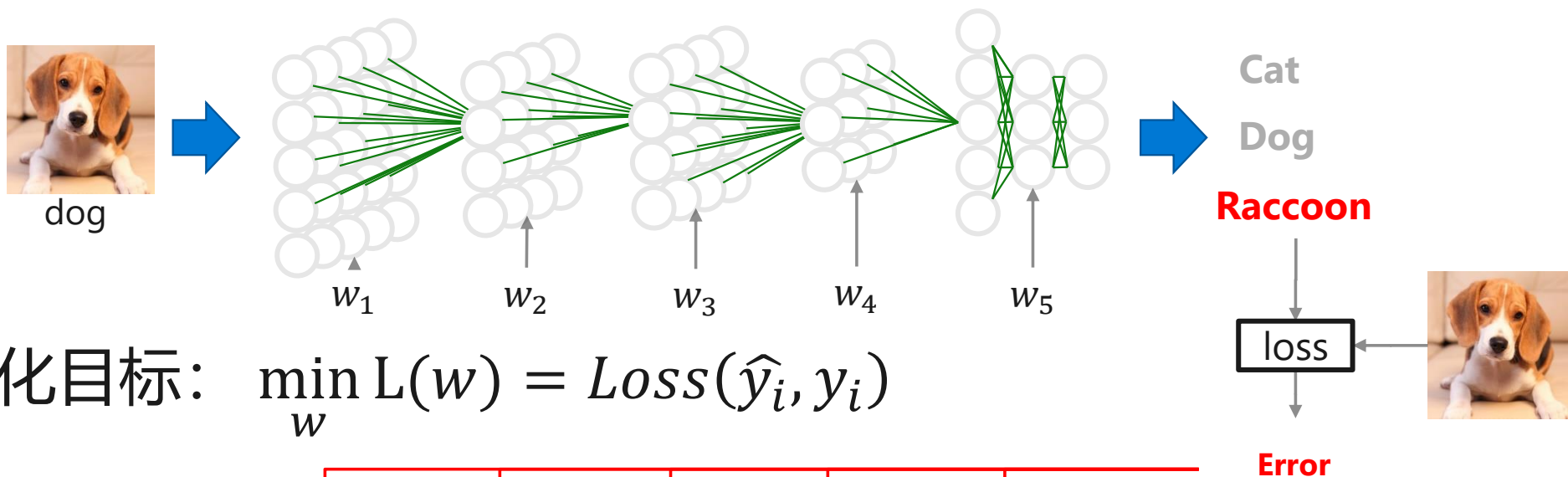
- Tensor
- Data-flow graph, DAG
- Backpropagation and auto-differentiation
- Graph execution and scheduling
- Symbolic and imperative execution, static vs dynamic graph
- Hardware device support

- 参考系统

- Caffe, Theano, DistBelief
- TensorFlow, CNTK
- PyTorch, Chainer, DyNet

# 回顾：深度学习基础

1. 定义一个带参数的函数（神经网络）： $\hat{y}_i = f(w, x_i)$



2. 定义优化目标： $\min_w L(w) = \text{Loss}(\hat{y}_i, y_i)$

$$\frac{d\text{error}}{dw_1} \quad \frac{d\text{error}}{dw_2} \quad \frac{d\text{error}}{dw_3} \quad \frac{d\text{error}}{dw_4} \quad \frac{d\text{error}}{dw_5}$$

3. 计算梯度并更新参数： $w_i \leftarrow w_i - \eta \nabla_{w_i} L(w)$

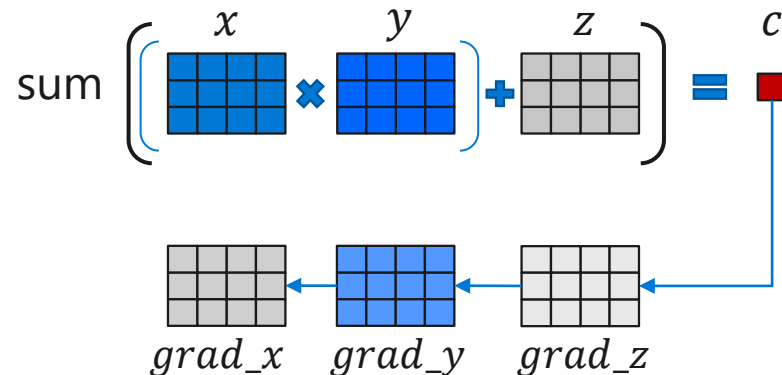
# 例：一种极端的计算方法

- 用某种高级语言从头实现一个模型的计算过程

```
import numpy as np
N, D = 3, 4
x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)

a = x * y
b = a + z
c = np.sum(b)

grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```



Efficiency



Python-like

Flexibility

# 例：另一种极端的计算方法

- 为常用模型在某个加速设备上实现一个高度优化的计算库

```
import xxlib  
  
x, y = load_data()  
  
y = xxlib.resnet152(x)
```



**Efficiency**



library

Python-like

**Flexibility**

# 深度学习计算框架的目的

- 提供**灵活**的编程模型和编程接口
  - 简洁的神经网络计算原语编程语言
  - 提供直观地模型构建方式
  - 较好的支持与现有生态环境融合
- 提供**高效**和可扩展的计算能力
  - 自动推导计算图
  - 自动编译优化算法，包括不限于：公共子表达式消除，内核融合，内存布局优化等
  - 根据不同体系结构和硬件设备自动并行化
  - 自动分布式化，并扩展到多个计算节点
  - 持续优化

# 早期的深度学习框架 (-2010)

## • 主要解决的问题

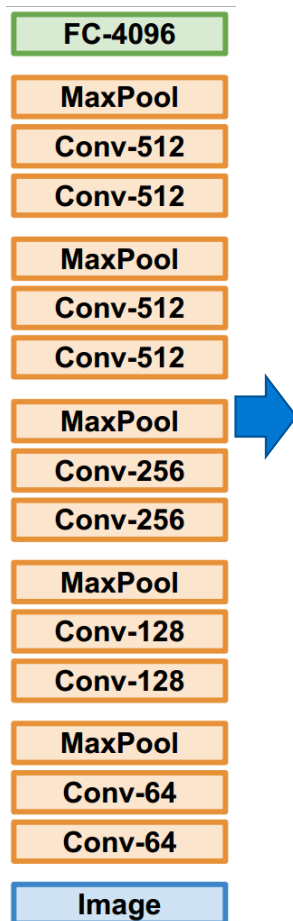
- 基于CNN网络的图像识别类模型为主，由一些常用的layers组成，如：
  - Convolution, Pooling, BatchNorm, Activation等

## • 主要特点

- 通过简单配置文件的形式定义神经网络
- 模型可由一些常用layer构成一个简单的图
- 框架提供每一个layer及其梯度计算实现
- 支持多设备加速：CPU和GPU的高效计算
- 代表框架：Caffe

## • 优点

- 提供了一定程度的可编程性
- 计算性能高：支持GPU加速计算



```
1 name: "ResNet-152"
2 input: "data"
3 input_dim: 1
4 input_dim: 3
5 input_dim: 224
6 input_dim: 224
7
8 layer {
9     bottom: "data"
10    top: "conv1"
11    name: "conv1"
12    type: "Convolution"
13    convolution_param {
14        num_output: 64
15        kernel_size: 7
16        pad: 3
17        stride: 2
18        bias_term: false
19    }
20 }
21
22 layer {
23     bottom: "conv1"
24     top: "conv1"
25     name: "bn_conv1"
26     type: "BatchNorm"
27     batch_norm_param {
28         use_global_stats: true
29     }
30 }
31
```

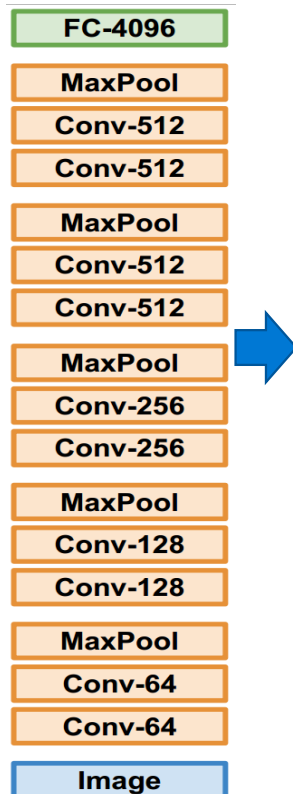
# 早期的深度学习框架 (-2010)

## • 主要特点

- 通过简单配置文件的形式定义神经网络
- 模型可由一些常用layer构成一个简单的图
- 框架提供每一个layer及其梯度计算实现
- 支持多设备加速：CPU和GPU的高效计算
- 代表框架：Caffe

## • 优点

- 提供了一定程度的可编程性
- 计算性能高：支持GPU加速计算



```
1 name: "ResNet-152"  
2 input: "data"  
3 input_dim: 1  
4 input_dim: 3  
5 input_dim: 224  
6 input_dim: 224  
7  
8 layer {  
9     bottom: "data"  
10    top: "conv1"  
11    name: "conv1"  
12    type: "Convolution"  
13    convolution_param {  
14        num_output: 64  
15        kernel_size: 7  
16        pad: 3  
17        stride: 2  
18        bias_term: false  
19    }  
20 }  
21  
22 layer {  
23     bottom: "conv1"  
24     top: "conv1"  
25     name: "bn_conv1"  
26     type: "BatchNorm"  
27     batch_norm_param {  
28         use_global_stats: true  
29     }  
30 }  
31
```

Efficiency

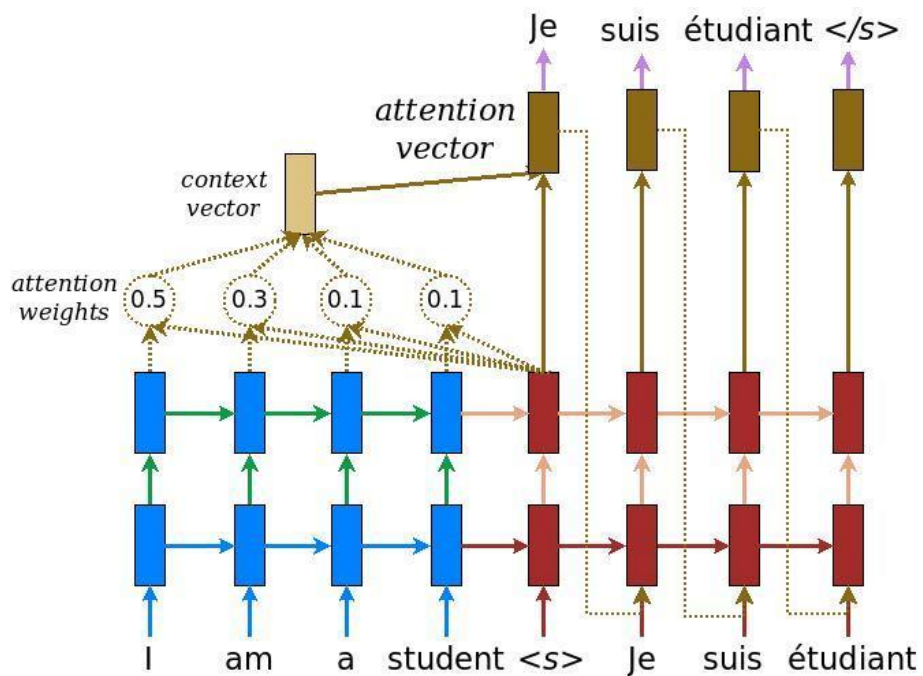


Flexibility



# 第一代框架的局限性 (I)

- 灵活性的限制难以满足深度学习的快速发展
  - 层出不穷的新型网络结构 (Layers) 要求针对每种Layer都要重新实现其前向和后向计算函数
  - 如attention layer, Batch normalization layer, sampled softmax等



```
Class AttentionLayer<CPU>
{
  void forward(inputs...)
  {
  }
  void backward(inputs, grad)
  {
  }
  ...
};
Class AttentionLayer<GPU>
{
  ...
};

REGISTER_LAYER("Attention", AttentionLayer);
```

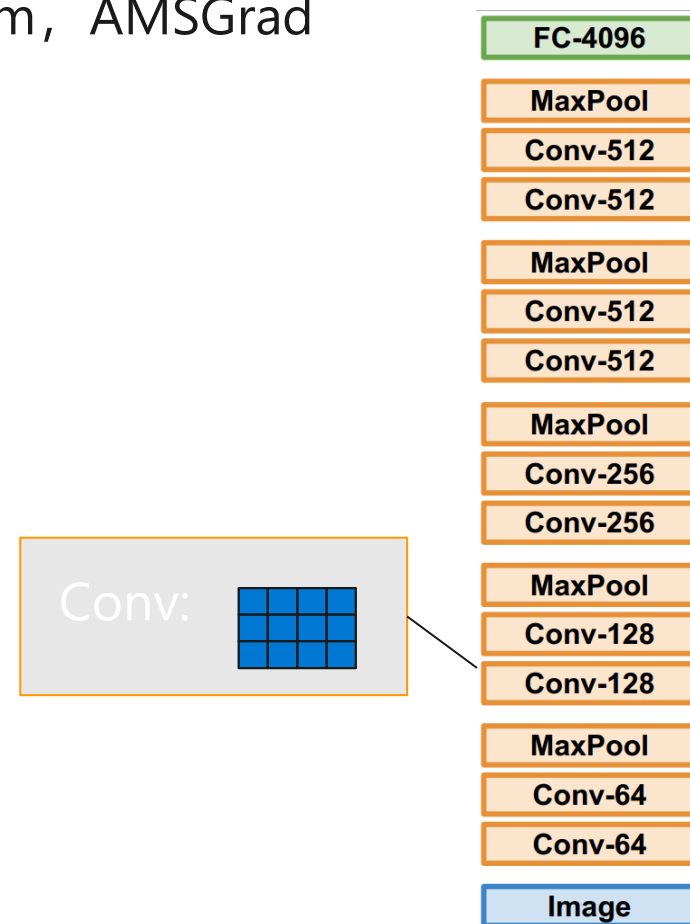
# 第一代框架的局限性(II)

- 新的优化器 (Optimizer) 要求对梯度和参数进行更通用复杂的运算
  - 如Adagrad, Adadelta, RMSprop, Adam, AdaMax, Nadam, AMSGrad

$$\text{SGD: } w \leftarrow w - \eta \nabla_w$$



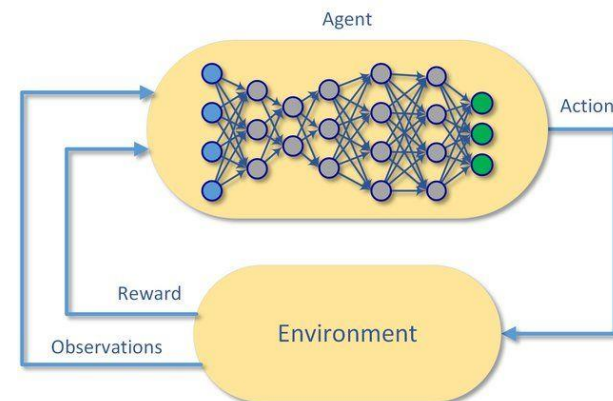
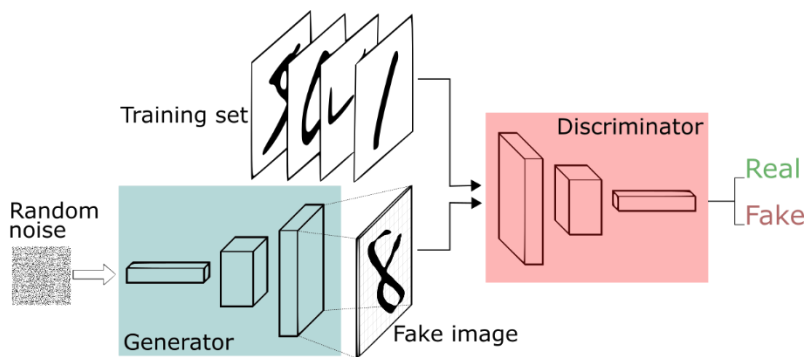
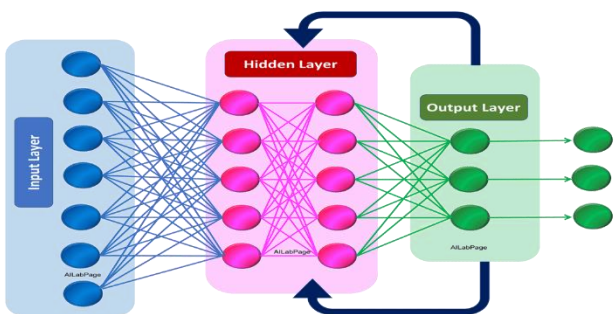
$$\text{SGD with momentum: } w \leftarrow w - (\gamma \nabla_w^{t-1} + \eta \nabla_w^t)$$



# 第一代框架的局限性(III)

- 基于简单的“前向+后向”的训练模式难以满足新的训练模式
  - 循环神经网络需要引入控制流，如RNN
  - 对抗神经网络需要两个网络交替训练
  - 强化学习模型需要和外部环境进行交互，如AlphaGo

## Recurrent Neural Networks



# 第二代深度学习框架 (2010-)

- 兼顾编程的灵活性和计算的高效性

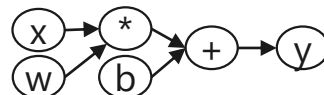


前端编程语言和接口

Python, Lua, R, C++

自动求导 (Auto Differentiation)

统一模型表示: 计算流图



图的优化与调度执行

Batching, Cache, Overlap

内核代码优化与编译

GPU kernel, auto kernel generation

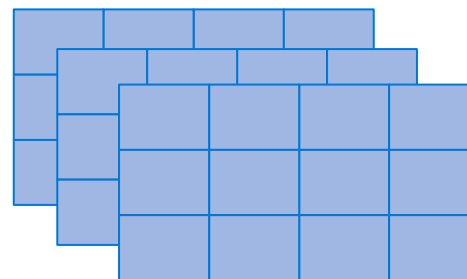
计算硬件

CPU, GPU, RDMA devices

# 基于数据流图 (DAG) 的计算框架

- 基本数据结构: Tensor (N维数组)

- Tensor形状: [2, 3, 4]
- 元素类型: int, float, string, etc.



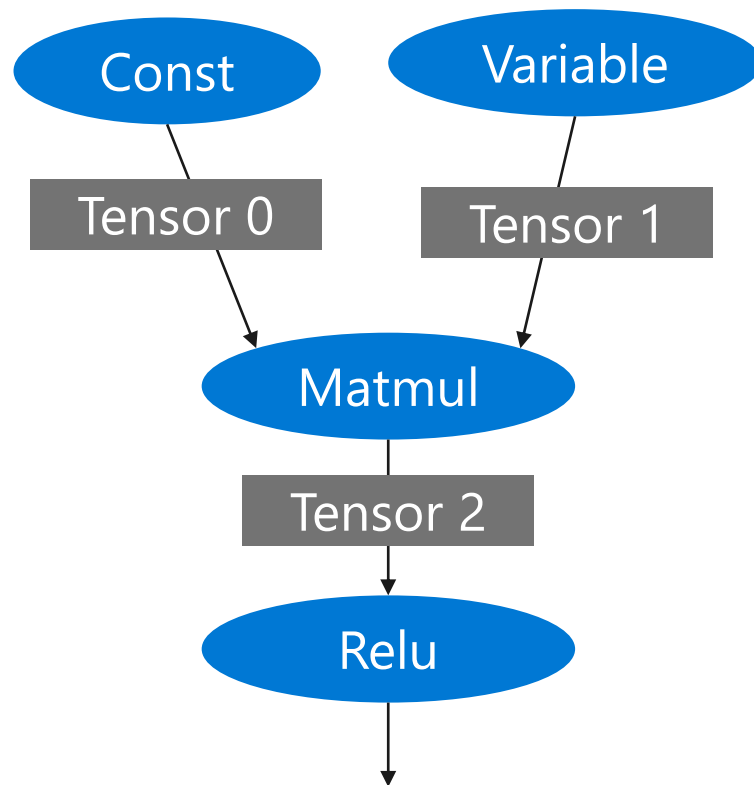
- 基本运算单元: Operator

- 由最基本的代数算子组成
- 每个Operator接收N个输入Tensor, 并输出M个输出Tensor
- TensorFlow中有>400个基本operator

Add	Log	While
Sub	MatMul	Merge
Mul	Conv	BroadCast
Div	BatchNorm	Reduce
Relu	Loss	Map
Tanh	Transpose	Reshape
Exp	Concatenate	Select
Floor	Sigmoid	....

# 基于数据流图 (DAG) 的计算框架

- 用数据流图表示计算逻辑和状态
  - 节点表示Operator
  - 边表示Tensor
- 计算状态 (如参数) 也是Operator
  - 如Variable Operator
- 特殊的Operator
  - 如: Switch, Merge, While 等用来构建控制流
- 特殊的边
  - 如: 控制边用来表示节点之间的依赖关系



# 编程语言与编程模型

## Numpy

```
import numpy as np
np.random.seed(0)
```

```
N, D = 3, 4
```

```
x = np.random.randn(N, D)
```

```
y = np.random.randn(N, D)
```

```
z = np.random.randn(N, D)
```

```
a = x * y
```

```
b = a + z
```

```
c = np.sum(b)
```

```
grad_c = 1.0
```

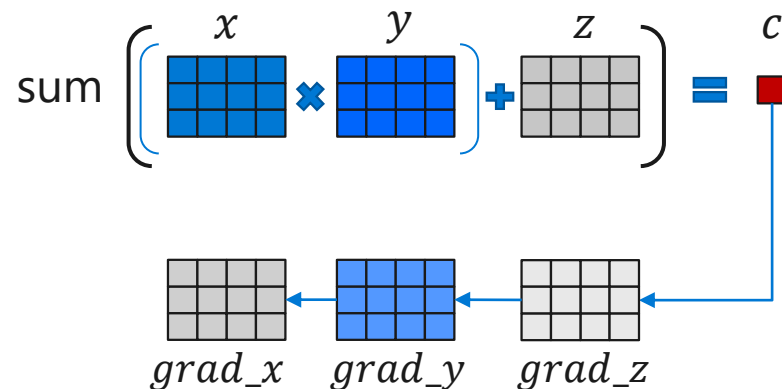
```
grad_b = grad_c * np.ones((N, D))
```

```
grad_a = grad_b.copy()
```

```
grad_z = grad_b.copy()
```

```
grad_x = grad_a * y
```

```
grad_y = grad_a * x
```



# 编程语言与编程模型

## Numpy

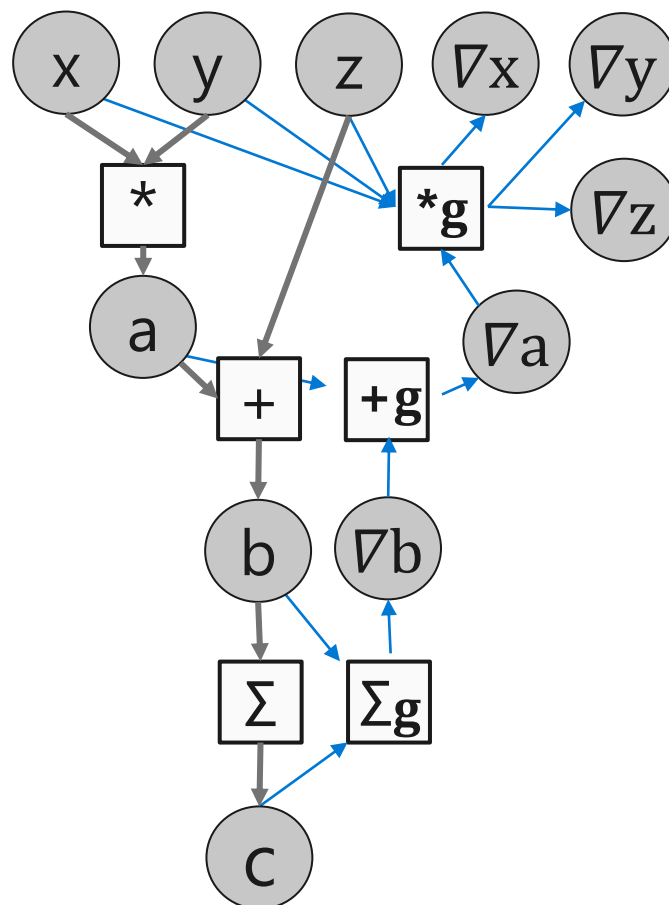
```
import numpy as np
np.random.seed(0)
```

```
N, D = 3, 4
```

```
x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)
```

```
a = x * y
b = a + z
c = np.sum(b)
```

```
grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```



## TensorFlow

```
import tensorflow as tf
np.random.seed(0)
```

```
N, D = 3, 4
```

```
x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)
```

```
a = x * y
b = a + z
c = tf.reduce_sum(b)
```

```
grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])
```

```
with tf.Session() as sess:
    sess.run([grad_z], feed_dict=values)
```



# 讨论:

- 数据流图表示深度学习模型的灵活性表现在哪些方面?
- 数据流图还有哪些好处?
- 你知道的其它基于数据流图的计算系统?

# 自动求导 (AD)

- 深度学习计算的核心—计算参数更新的梯度

$$L(w) = \text{Loss}(f(w, x_i), y_i) \rightarrow \frac{\partial L(w)}{\partial w}$$

- 求导计算是一个经典的问题

$$L(x) = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

# 符号求导(Symbolic Differentiation)

- 根据简单函数的导数公式，以及导数变换公式精确的计算出一个复杂函数的导数形式化表示

- 导数转换公式:  $\frac{d}{dx}(f(x) + g(x)) = \frac{d}{dx}f(x) + \frac{d}{dx}g(x)$

$$\frac{d}{dx}(f(x)g(x)) = \left(\frac{d}{dx}f(x)\right)g(x) + \left(\frac{d}{dx}g(x)\right)f(x)$$

- 例如:  $L(x) = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$ 的导数可以推导出为

- $\frac{dL}{dx} = \exp(\exp(x) + \exp(x)^2)((\exp(x) + 2\exp(x)^2) + \cos(\exp(\exp(x) + \exp(x)^2) (\exp(x) + 2\exp(x)^2)$

- 在深度学习中的应用问题

- 深度学习网络非常大 → 待求导函数复杂 → 难以高效的求解
- 深度中一些算子无法求导: 如Relu, Switch

# 数值求导(Numerical Differentiation)

- 通过数值逼近的方法计算近似导数

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x + he_i) - f(x)}{h}$$

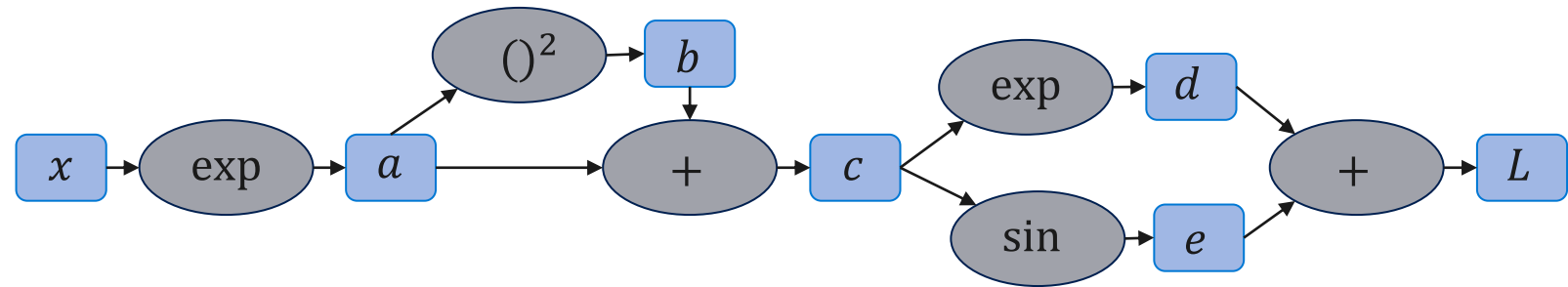
- 在深度学习中的应用问题
  - 由于数值计算中的截断和近似问题导致无法得到精确导数
  - 同样无法适用于深度中一些无法求导的算子：如Relu, Switch

# 自动求导 (Auto Differentiation)

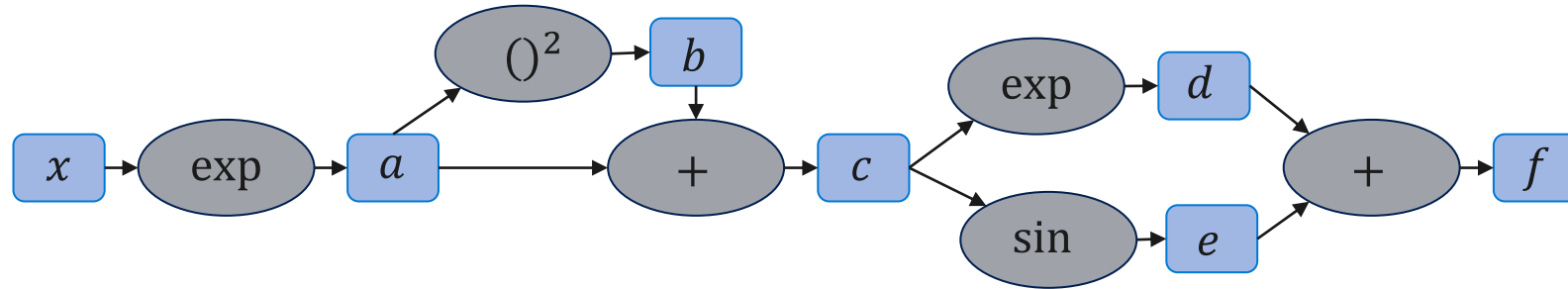
- 通过引入中间变量将一个复杂的函数分解成一系列基本函数
- 将这些基本函数构成一个计算流图

$$L(x) = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

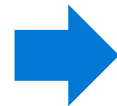
$$\begin{aligned} a &= \exp(x) \\ b &= a^2 \\ c &= a + b \\ d &= \exp(c) \\ e &= \sin(c) \\ L &= d + e \end{aligned}$$



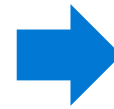
# 自动求导 (Auto Differentiation) (II)



$$\begin{aligned} a &= \exp(x) \\ b &= a^2 \\ c &= a + b \\ d &= \exp(c) \\ e &= \sin(c) \\ f &= d + e \end{aligned}$$

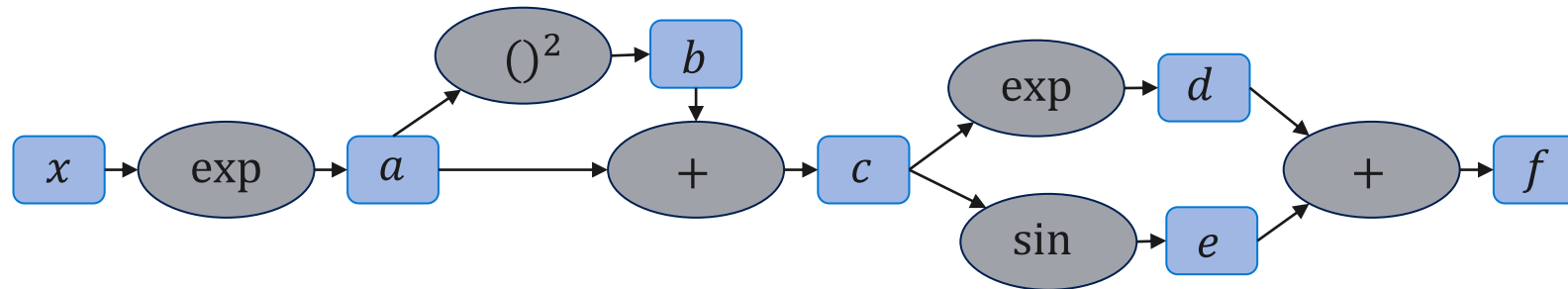


$$\begin{aligned} \frac{df}{dd} &= 1 \\ \frac{df}{de} &= 1 \\ \frac{df}{dc} &= \frac{df}{dd} \frac{dd}{dc} + \frac{df}{de} \frac{de}{dc} \\ \frac{df}{db} &= \frac{df}{dc} \frac{dc}{db} \\ \frac{df}{da} &= \frac{df}{dc} \frac{dc}{da} + \frac{df}{db} \frac{db}{da} \\ \frac{df}{dx} &= \frac{df}{da} \frac{da}{dx} \end{aligned}$$



$$\begin{aligned} \frac{df}{dd} &= 1 \\ \frac{df}{de} &= 1 \\ \frac{df}{dc} &= \frac{df}{dd} \exp(c) + \frac{df}{de} \cos(c) \\ \frac{df}{db} &= \frac{df}{dc} \\ \frac{df}{da} &= \frac{df}{dc} + \frac{df}{db} 2a \\ \frac{df}{dx} &= \frac{df}{da} \exp(x). \end{aligned}$$

# 自动求导 (AD) (III)



$$\begin{aligned} a &= \exp(x) \\ b &= a^2 \\ c &= a + b \\ d &= \exp(c) \\ e &= \sin(c) \\ f &= d + e \end{aligned}$$



$$\begin{aligned} \frac{df}{dd} &= 1 \\ \frac{df}{de} &= 1 \\ \frac{df}{dc} &= \frac{df}{dd} \exp(c) + \frac{df}{de} \cos(c) \\ \frac{df}{db} &= \frac{df}{dc} \\ \frac{df}{da} &= \frac{df}{dc} + \frac{df}{db} 2a \\ \frac{df}{dx} &= \frac{df}{da} \exp(x). \end{aligned}$$

Forward propagation

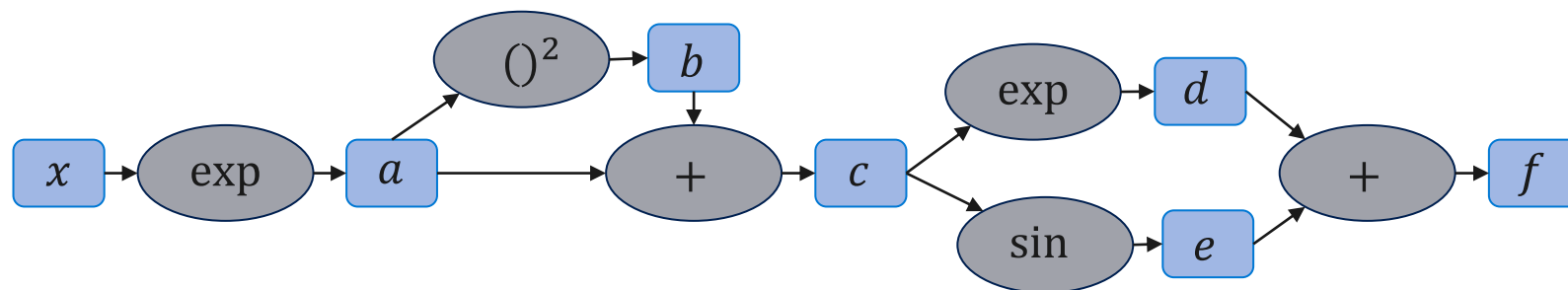
Backward propagation

# 讨论:

- Forward 和 Backward propagation 计算量是否一样?
- 深度学习中为什么大多使用 Backward propagation?



# 如何在深度学习框架中实现自动求导



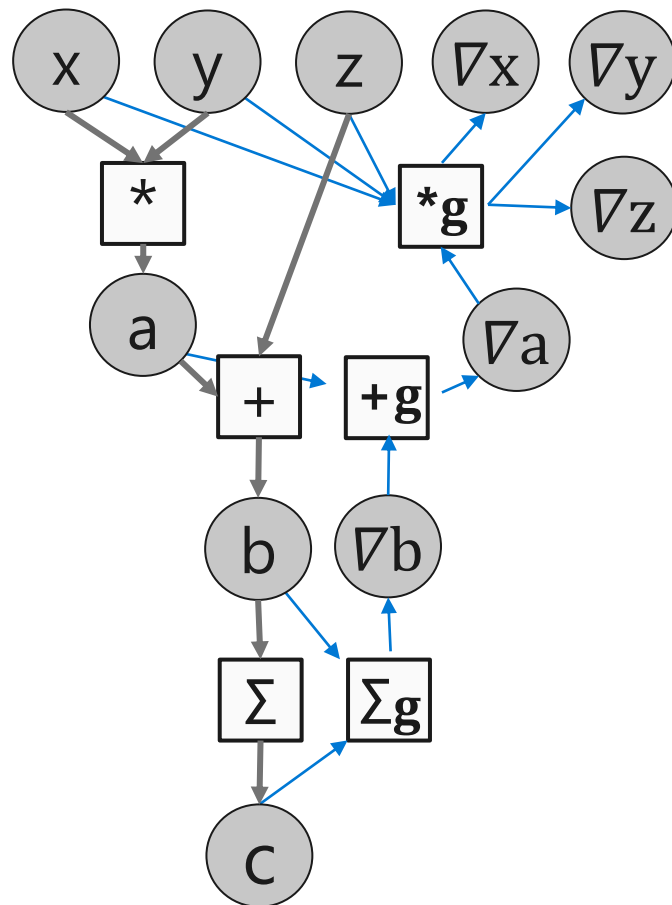
前向计算并保留中间计算结果  
根据BP的原理依次计算出中间导数

主要问题：  
需要保存大量中间计算结果

$$\begin{aligned} \frac{df}{dd} &= 1 \\ \frac{df}{de} &= 1 \\ \frac{df}{dc} &= \frac{df}{dd} \exp(c) + \frac{df}{de} \cos(c) \\ \frac{df}{db} &= \frac{df}{dc} \\ \frac{df}{da} &= \frac{df}{dc} + \frac{df}{db} 2a \\ \frac{df}{dx} &= \frac{df}{da} \exp(x). \end{aligned}$$

# 如何在深度学习框架中实现自动求导 (II)

将导数的计算也表示成数据流图



优点:  
方便全局图优化  
节省内存

## TensorFlow

```
import tensorflow as tf
np.random.seed(0)

N, D = 3, 4

x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)

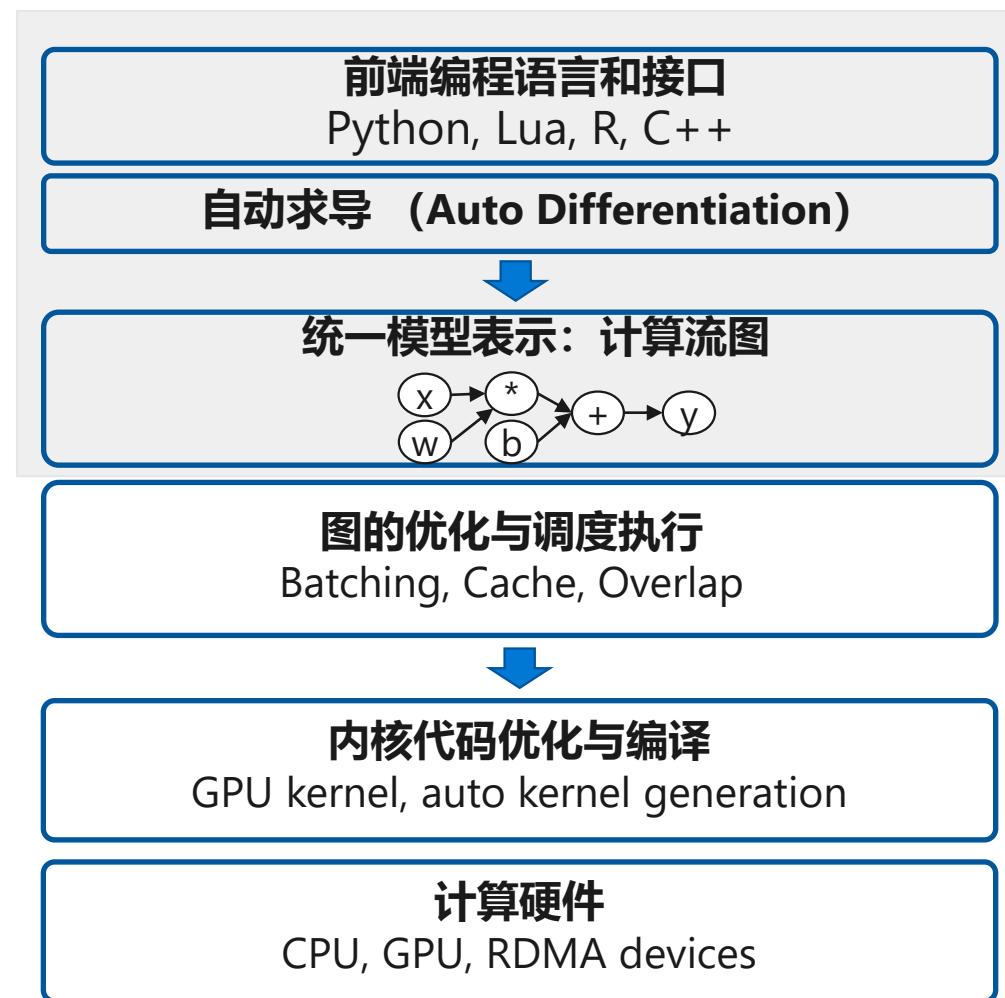
a = x * y
b = a + z
c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    sess.run([grad_z], feed_dict=values)
```

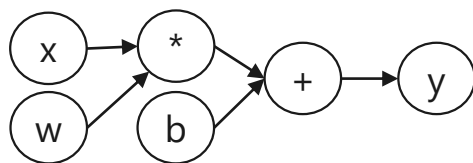
# 小结:

- **模型表示**: 数据流图
- **前端语言**: 用来构建数据流图
- **自动求导**: 基于backpropagation的原理自动构建求导数据流图



# 图优化

- “先定义后执行”的模式允许框架在计算前看到全图信息
- 数据流图作为深度学习框架中的高层中间表示，可以允许任何等价图优化Pass去化简计算流图或提高执行效率



表达式化简

公共子表达式消除

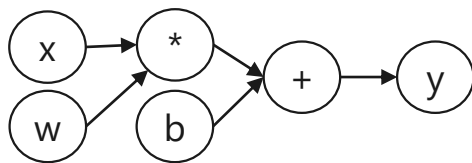
常数传播

Operator Batch

表达式替换

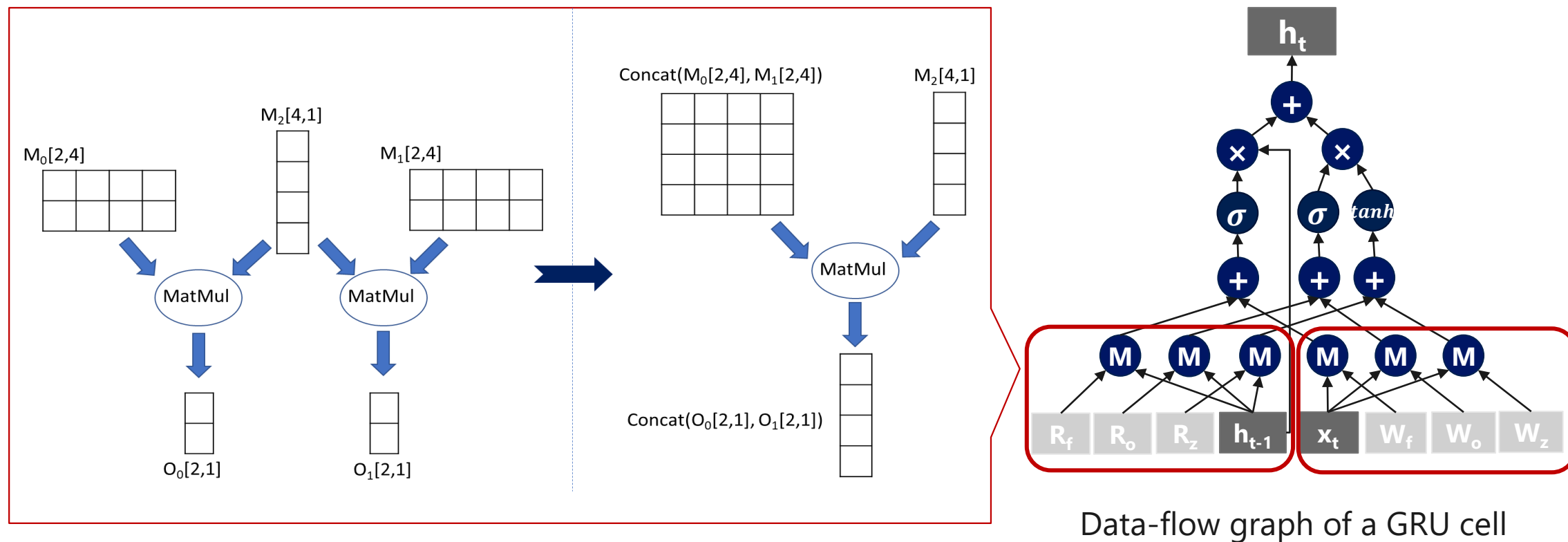
算子融合

...



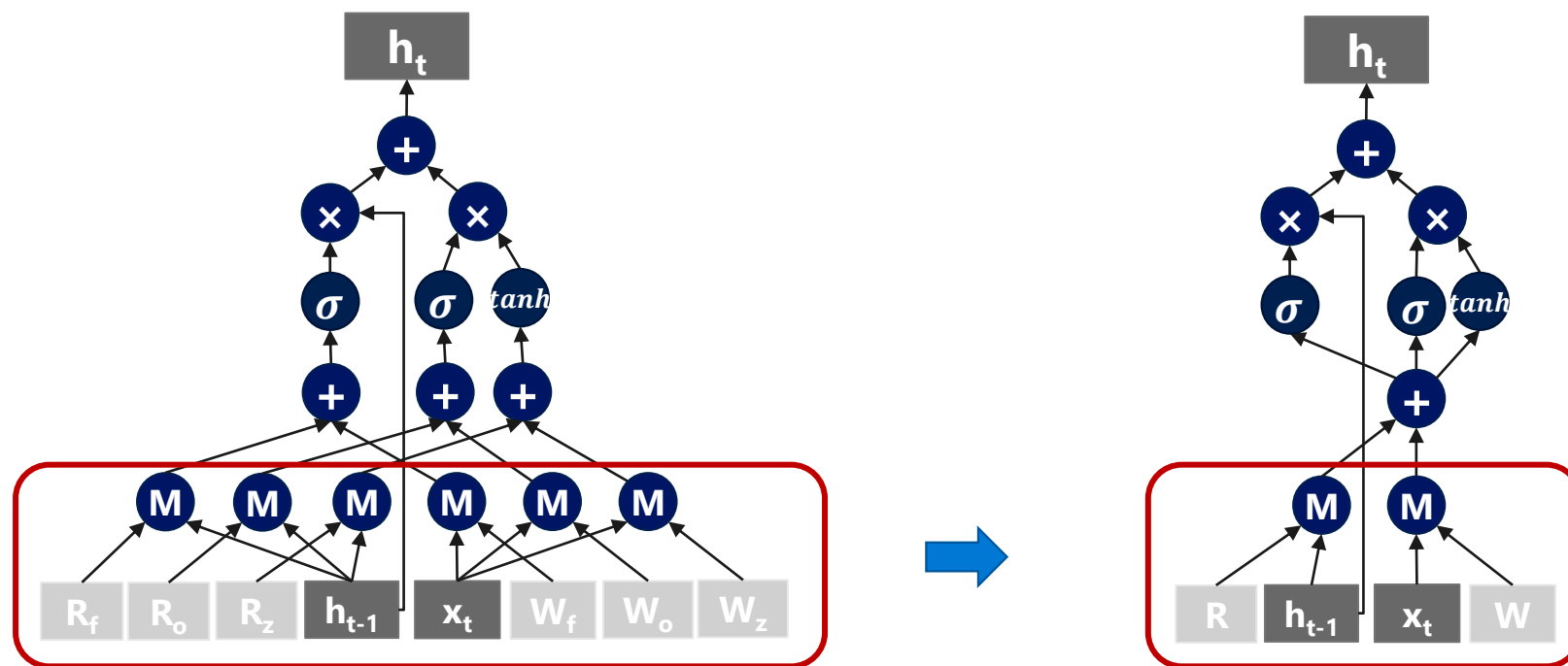
# 图优化：GEMM自动融合

- Batch same-type operators to leverage GPU massive parallelism



# 图优化：GEMM自动融合

- Batch same-type operators to leverage GPU massive parallelism

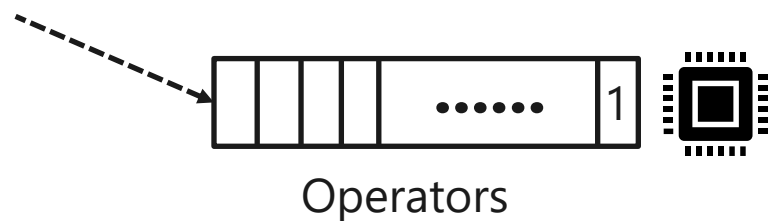
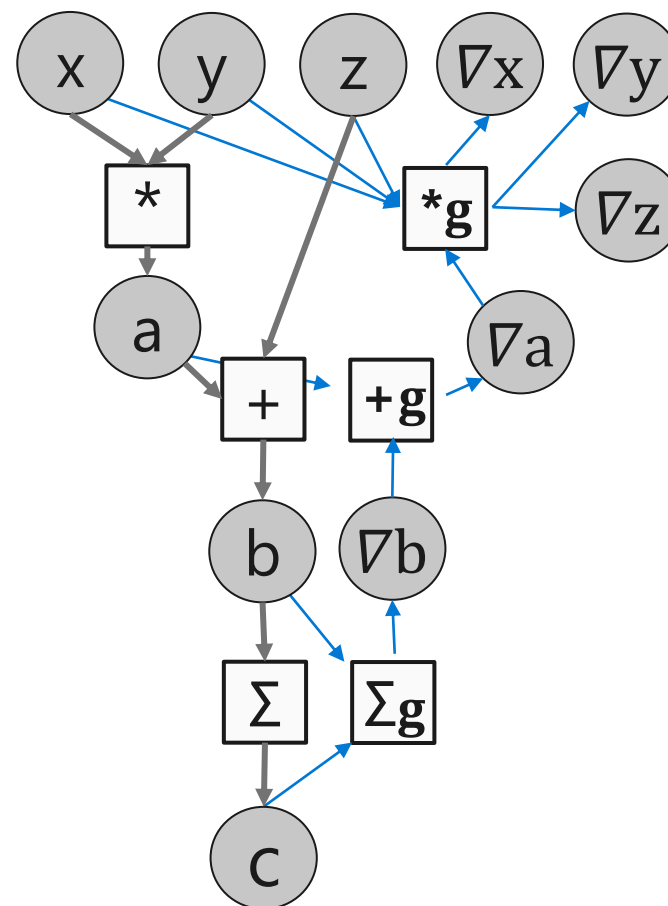


Data-flow graph of a GRU cell

# 数据流图的调度与执行

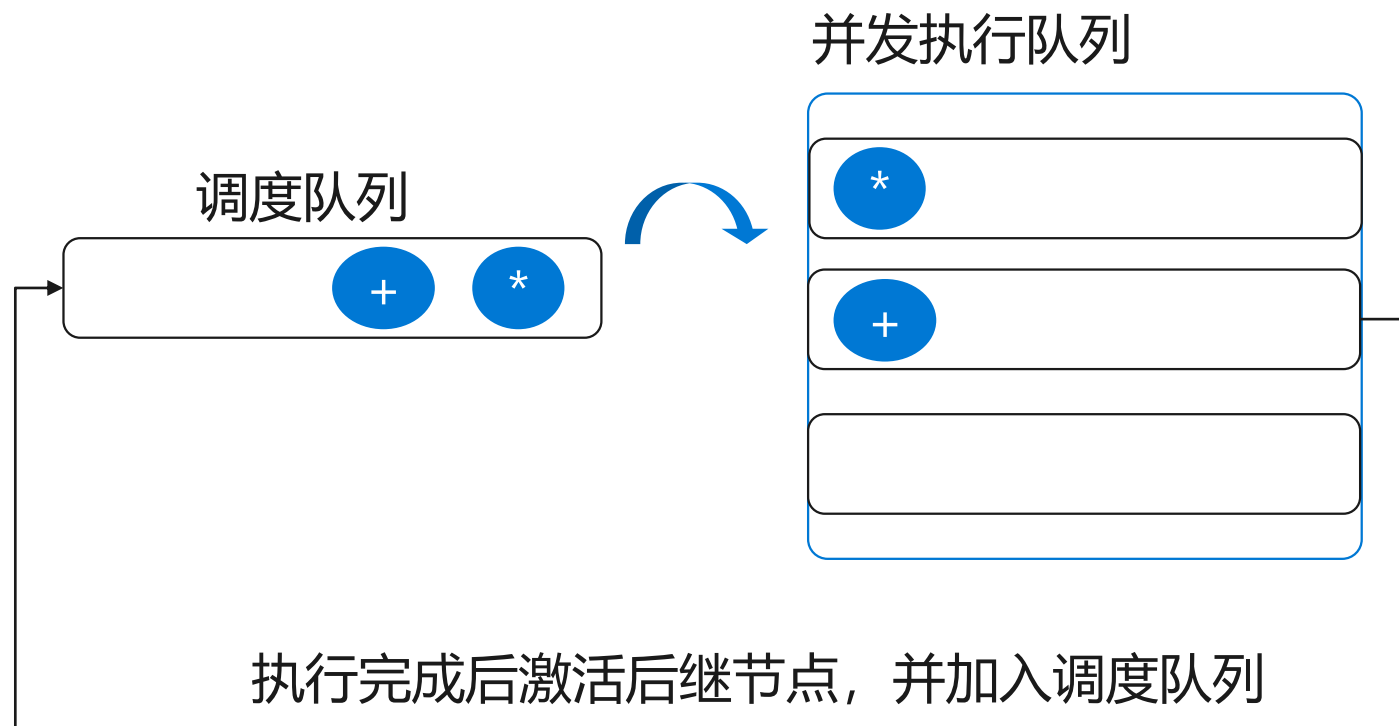
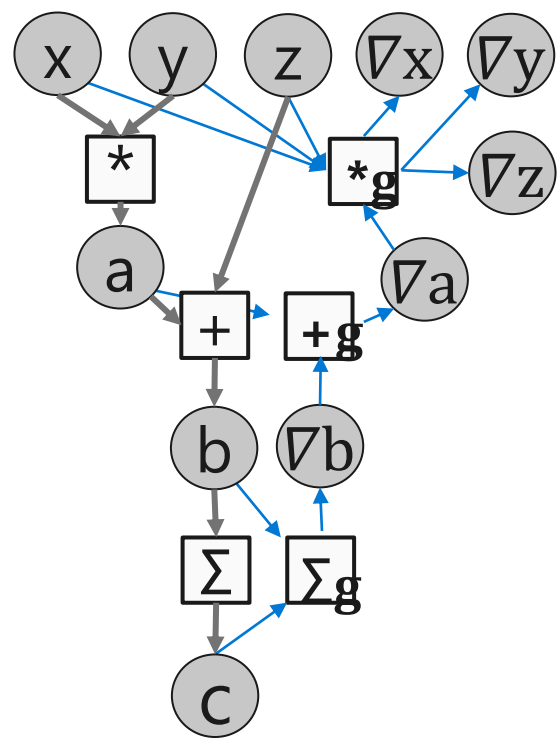
根据依赖关系，依次调度运行代码

1.  $\text{Multiply}(x, y) \rightarrow a$
2.  $\text{Add}(z, a) \rightarrow b$
3.  $\text{ReduceSum}(b) \rightarrow c$
4.  $\text{ReduceSum\_grad}(c, b) \rightarrow b\_delta$
5.  $\text{Sum\_grad}(b\_delta, a) \rightarrow a\_delta$
6.  $\text{Add\_grad}(a\_delta, z) \rightarrow z\_delta$
7.  $\text{Multiply}(a\_delta, x, y) \rightarrow x\_delta, y\_delta$



# 数据流图的并发执行

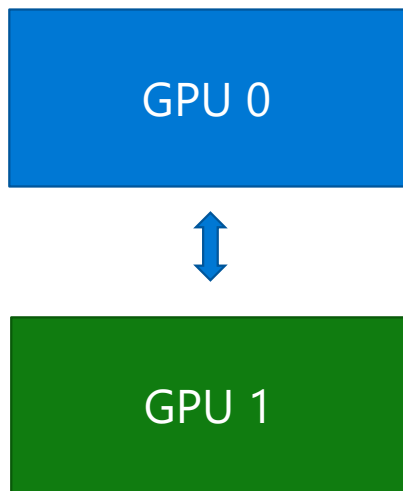
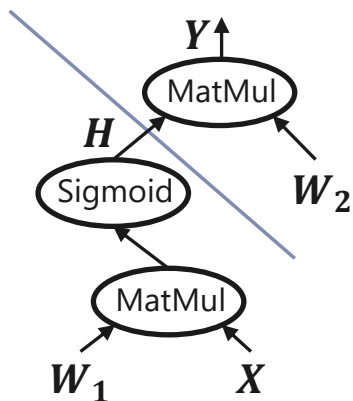
- 数据流图准确的描述了算子之间的依赖关系
- 根据数据流图找到相互独立的算子进行并发调度，提高计算的并行性





# 数据流图的划分与设备放置

- 显式图划分



```
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3000, 4000

with tf.device('/gpu:0'):
    x = tf.placeholder(tf.float32)
    y = tf.placeholder(tf.float32)
    z = tf.placeholder(tf.float32)

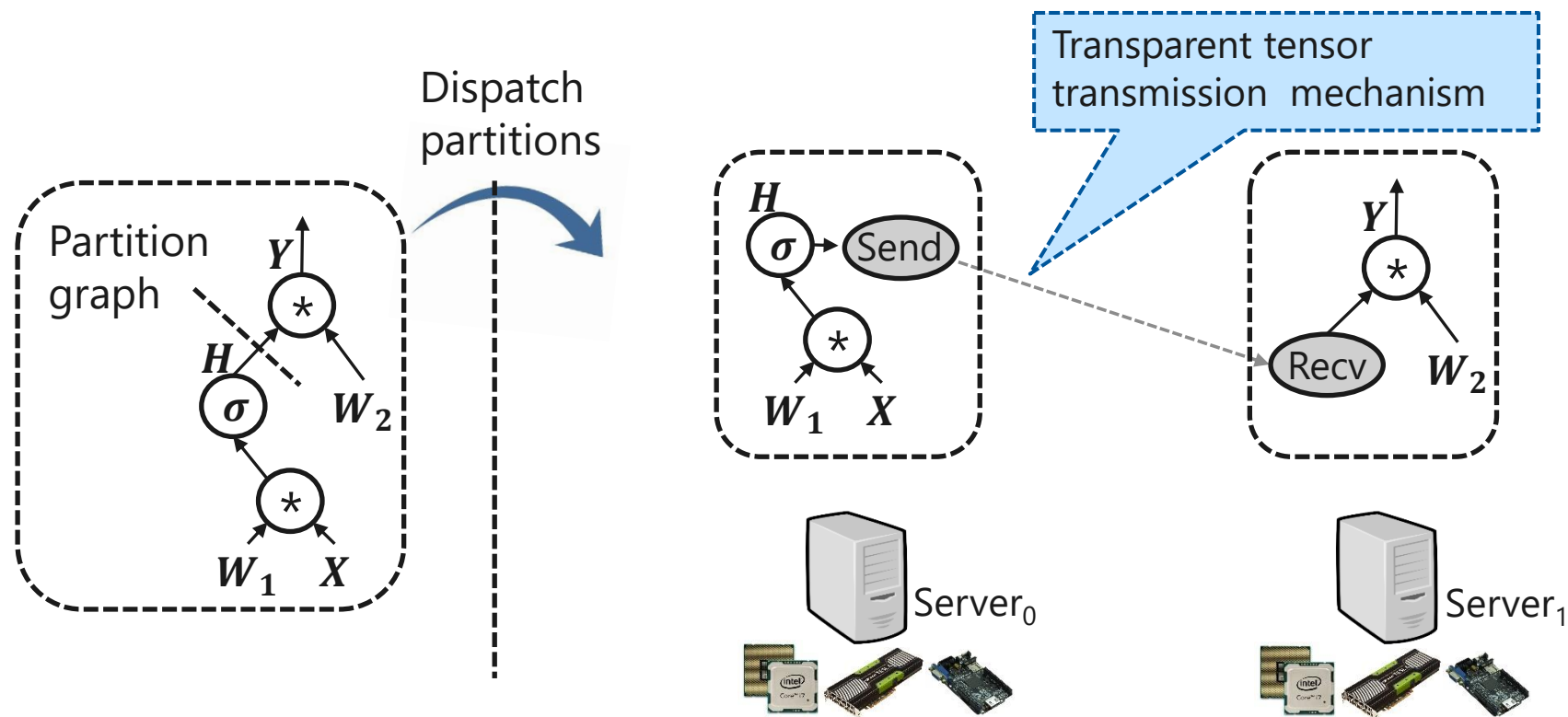
    a = x * y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                    feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

# 数据流图的划分与设备放置

- 跨设备的边将被自动替换成一组Send/Recv operators



# 计算内核(Kernel)与多硬件支持

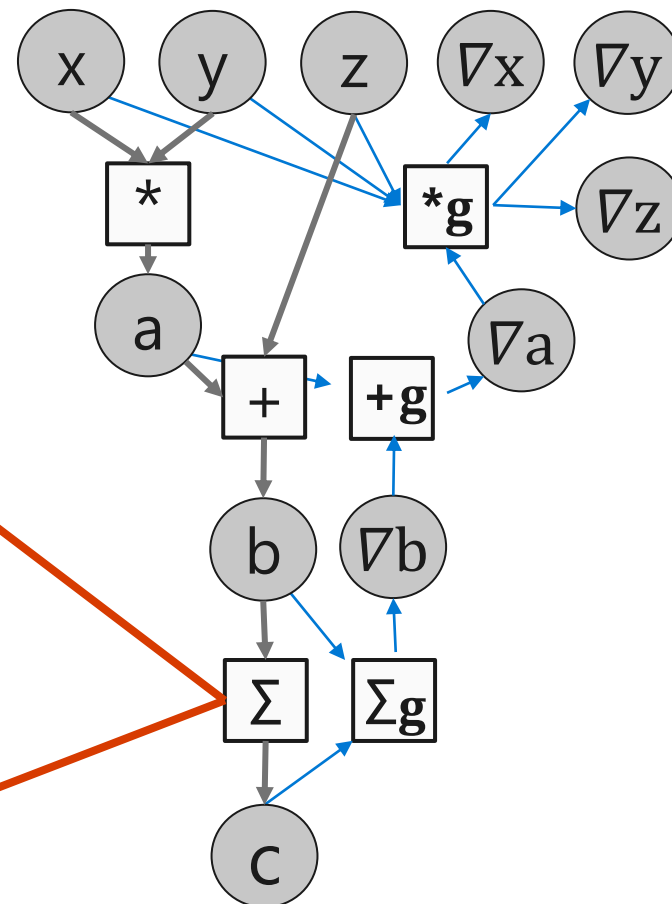
- 内核(Kernel)是定义了一个算子在某种具体设备的计算实现
- 每个Operator都可以注册多个Kernel
  - 根据计算设备不同可以有：GPU kernel, CPU kernel, ASIC kernel, FPGA kernel
  - 根据数据类型不同可以有：float kernel, int kernel, half kernel等
  - 根据Operator属性不同也可以有多种kernel
- 运行时框架会自动根据Operator的设备类型、数据类型和属性选择对应的kernel来执行

# 例：TensorFlow中注册的CPU/GPU kernel

```

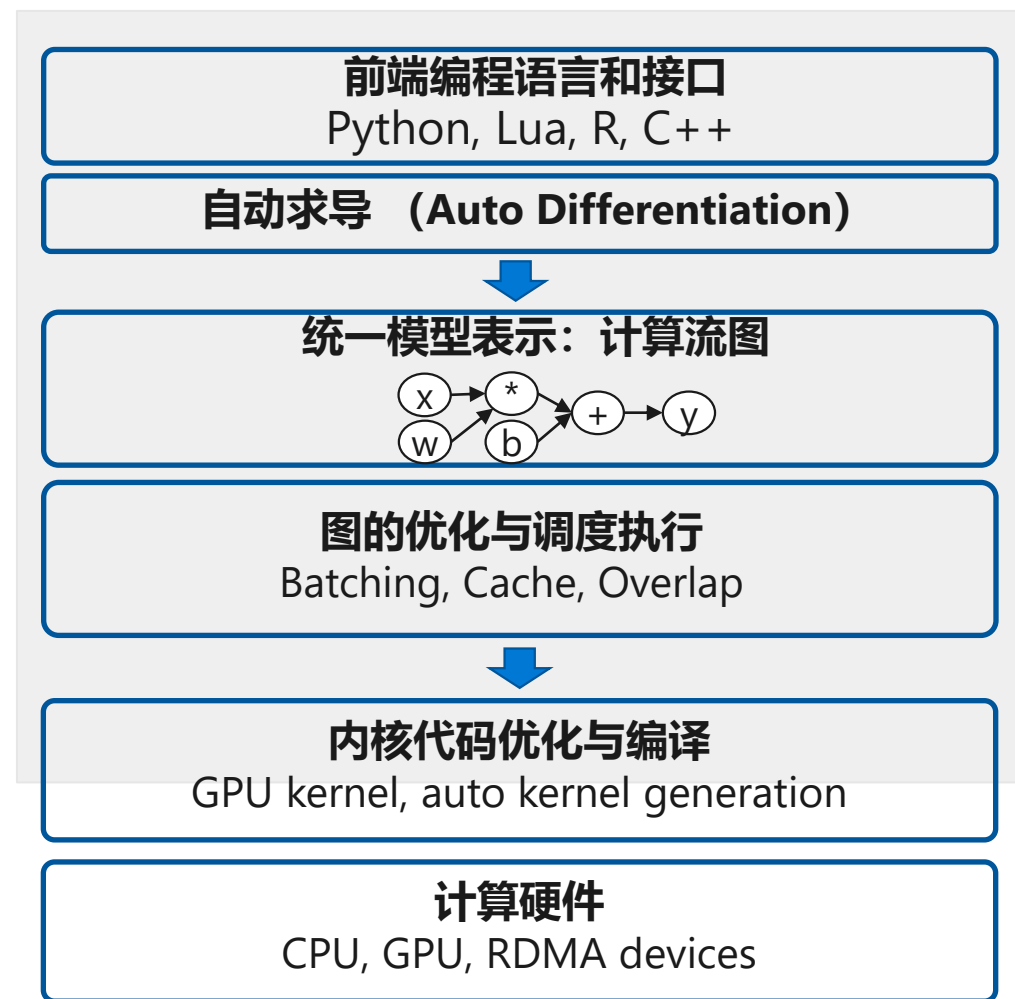
1  template <typename Device, typename OUT_T, typename IN_T,
2      typename ReductionAxes, typename Scalar>
3  struct ReduceEigenImpl<Device, OUT_T, IN_T, ReductionAxes,
4      functor::MeanReducer<Scalar>> {
5      void operator()(const Device& d, OUT_T out, IN_T in,
6          const ReductionAxes& reduction_axes,
7          const functor::MeanReducer<Scalar>& reducer) {
8          static_assert(std::is_same<Scalar, typename OUT_T::Scalar>::value, "");
9          Eigen::internal::SumReducer<Scalar> sum_reducer;
10         out.device(d) = in.reduce(reduction_axes, sum_reducer) / CPU code
11         static_cast<Scalar>(in.size() / out.size());
12     }
13 };
14
15 // T: the data type
16 // REDUCER: the reducer functor
17 // NUM_AXES: the number of axes to reduce
18 // IN_DIMS: the number of dimensions of the input tensor
19 #define DEFINE(T, REDUCER, IN_DIMS, NUM_AXES)
20     template void ReduceFunctor<GPUDevice, REDUCER>::Reduce( GPU code
21         OpKernelContext* ctx, TTypes<T, IN_DIMS - NUM_AXES>::Tensor out, \
22         TTypes<T, IN_DIMS>::ConstTensor in, \
23         const Eigen::array<Index, NUM_AXES>& reduction_axes, \
24         const REDUCER& reducer);

```



# 小结:

- **模型表示:** 数据流图
- **前端语言:** 用来构建数据流图
- **自动求导:** 基于backpropagation的原理自动构建求导数据流图
  
- **图的优化:** 图化简
- **调试执行:** 并发调度
- **设备放置:** 显式和隐式图切分
- **算子内核:** 模块化支持多设备



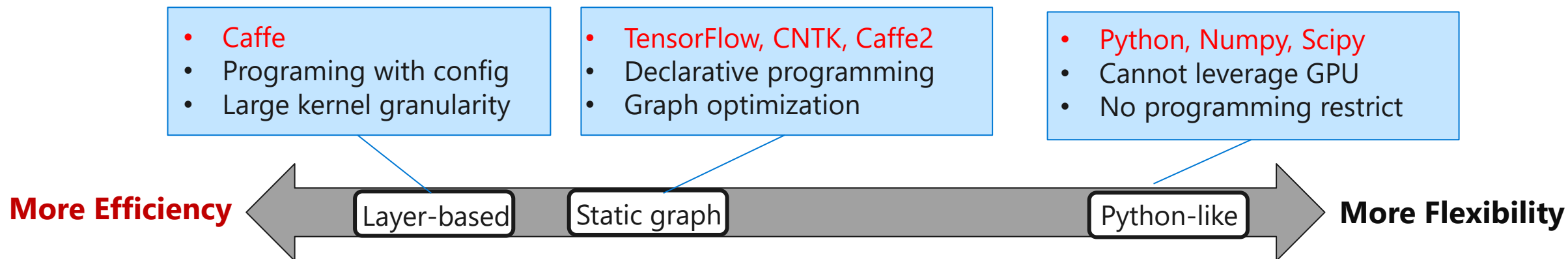
# 基于静态数据流图的计算框架的优点和缺点

## · 优点

- 计算效率较高：静态定义图后可以全局图优化
- 内存使用效率高：可以根据数据流图准确分析内存使用生命周期
- 较好的灵活性：可以表示大部分可由基本算子组成的网络

## · 缺点

- 可调试性：先申明后执行导致不能实时得到计算结果，难以Debug
- 表达的灵活性：受限于算子集合，如对于带有控制流的网络支持不友好或无法支持



# 基于动态数据流图的计算框架

- 边定义边执行-Define-by-run
  - 不预先定义计算流图
- 特点：
  - 模型可由任意Layer构成
  - 用户自己定义Layer可以内置layer组成
  - 可实时得到计算结果
  - 控制流由高层语言表示：如Python
  - 支持多设备加速：CPU和GPU的高效计算
  - 代表框架：PyTorch
- 优点
  - 提供了灵活的可编程性和可调试性

## PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D).cuda())
y = Variable(torch.randn(N, D).cuda())
z = Variable(torch.randn(N, D).cuda())

for i in range(10):
    a = x * y
    b = a + z
    c = c + torch.sum(b)

c.backward()
```

# 基于动态数据流图的计算框架： PyTorch Example

```
class LinearLayer(Module):  
    def __init__(self, in_sz, out_sz):  
        super().__init__()  
        t1 = torch.randn(in_sz, out_sz)  
        self.w = nn.Parameter(t1)  
        t2 = torch.randn(out_sz)  
        self.b = nn.Parameter(t2)  
  
    def forward(self, activations):  
        t = torch.mm(activations, self.w)  
        return t + self.b
```

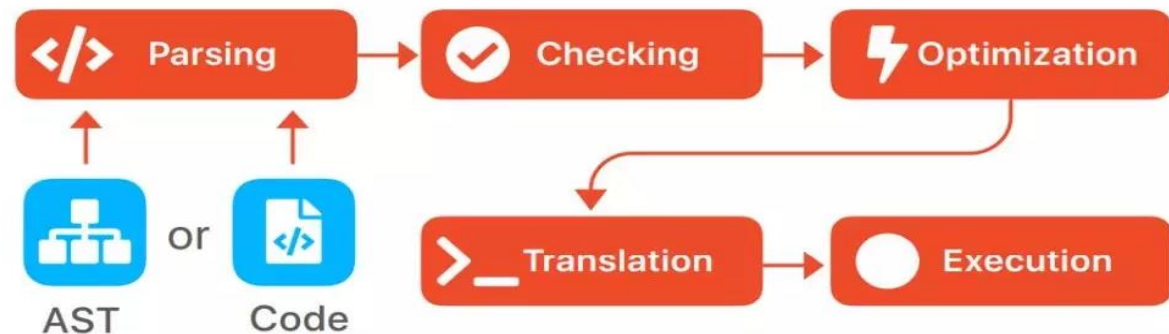
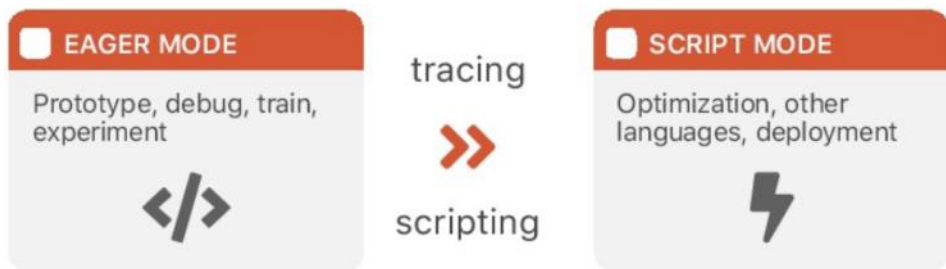
```
class FullBasicModel(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv = nn.Conv2d(1, 128, 3)  
        self.fc = LinearLayer(128, 10)  
  
    def forward(self, x):  
        t1 = self.conv(x)  
        t2 = nn.functional.relu(t1)  
        t3 = self.fc(t2)  
        return nn.functional.softmax(t3)
```



# 执行性能存在的问题与优化

- 由于丢失了全局数据流图，天然损失了全局图优化的好处
  - 无法进行全图化简优化
  - 无法进行全图精确内存分配管理
- 先天不足、后天来补：
  - **高效的C++核心**
    - 所以核心Layer的后端都是通过C++实现，对GPU也有较好的内核实现
  - **分离控制流和数据流的执行**
    - 控制流通过高级Python语言来执行，计算量较大的数据流部分再由C++执行
    - 通过GPU的异步接口调度GPU内核，从而实现调度和执行逻辑的重叠
  - **Tensor分配的缓存机制**
    - 通过对已分配的Tensor进行缓存来减少动态内存分配时的开销
  - **多进程执行**
    - 使用多进程来避免Python多线程的同步锁

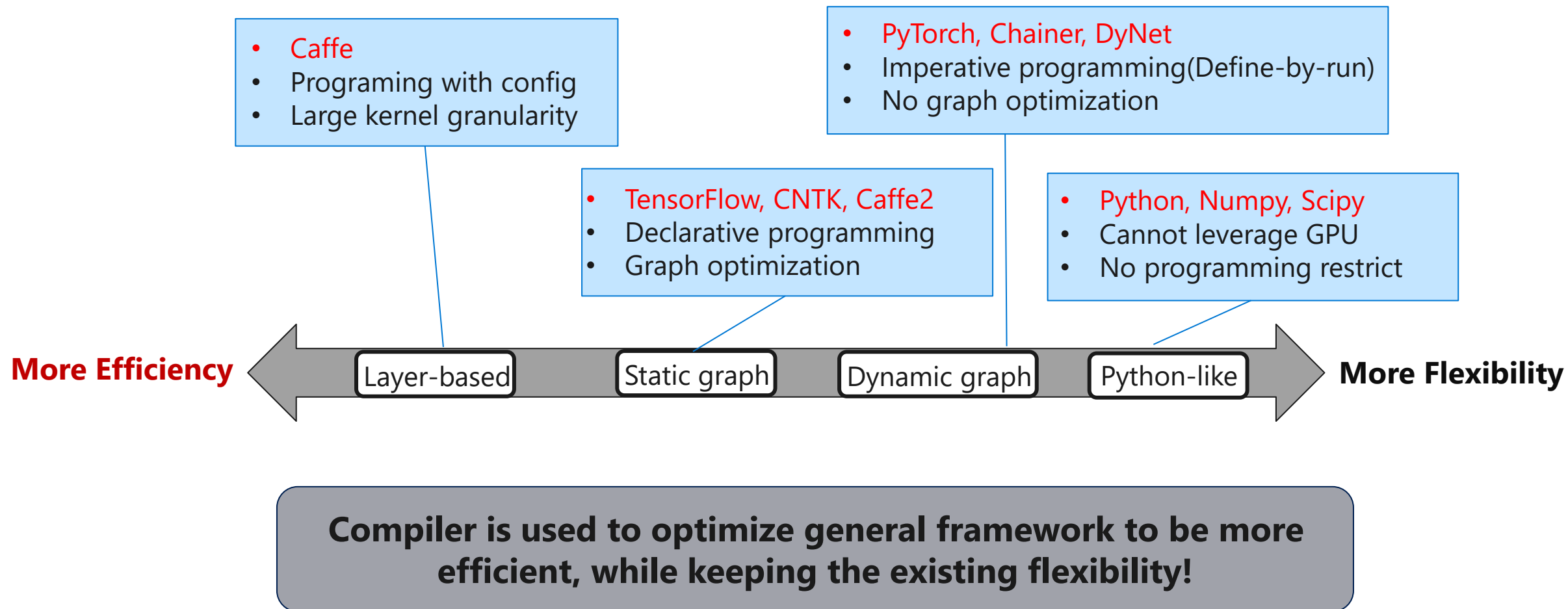
# 执行效率优化-JIT



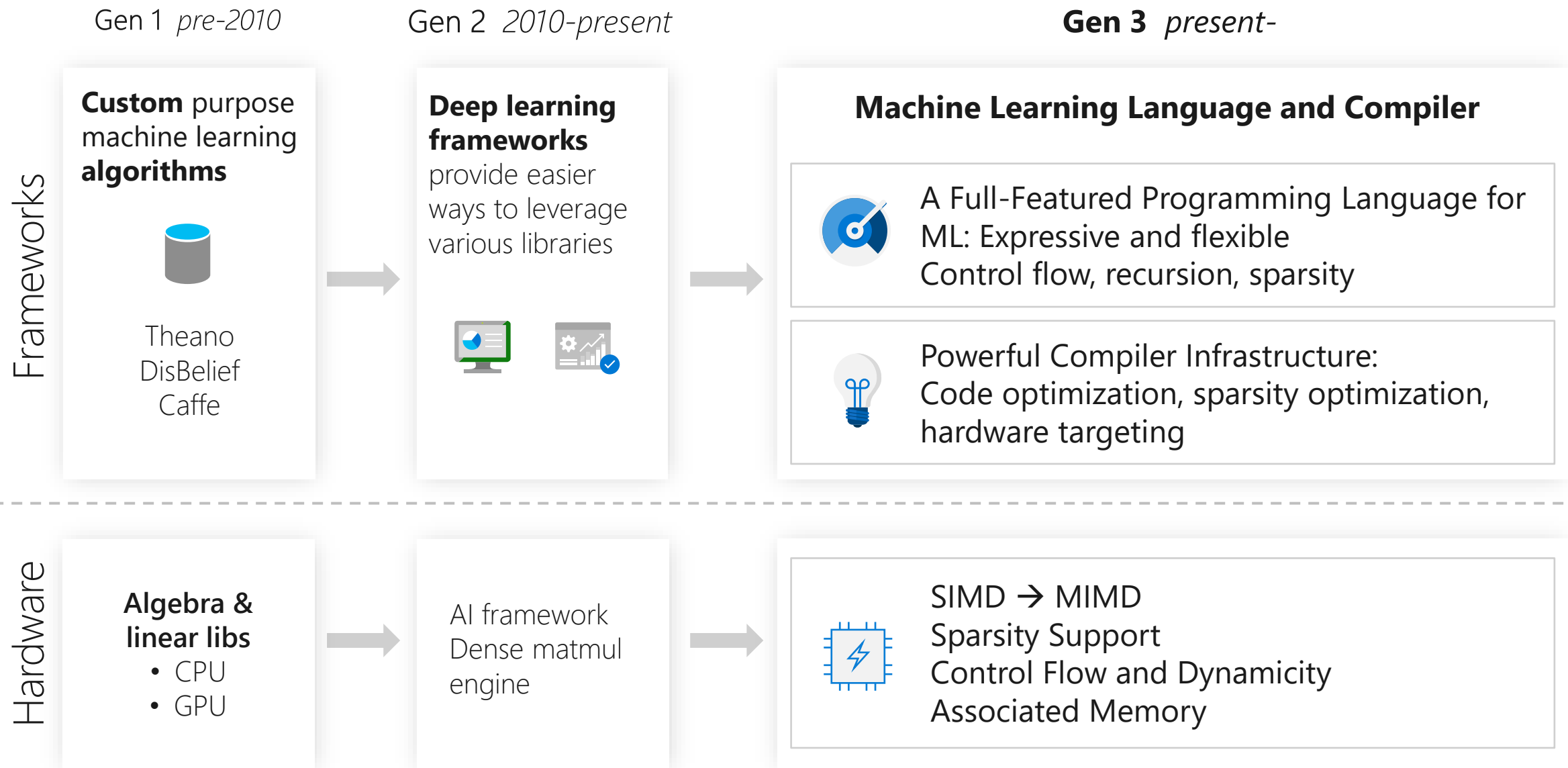
```
def dumb_function(x):  
    return x.t().t()
```

```
>>> traced_fn = torch.jit.trace(dumb_function,  
...                             torch.ones(2,2))  
>>> traced_fn.graph_for(torch.ones(2,2))  
graph(%x : Float(*, *)) {  
    return (%x);  
}
```

# 小结



# 计算框架的进步



# 参考阅读

- Large Scale Distributed Deep Networks, NIPS'12
- Caffe: Convolutional Architecture for Fast Feature Embedding, MM'14
- TensorFlow: A System for Large-Scale Machine Learning, OSDI'16
- PyTorch: An Imperative Style, High-Performance Deep Learning Library, NIPS'19
- Theano: A Python framework for fast computation of mathematical expressions, arXiv 16
- Automatic Differentiation in Machine Learning: a Survey, [Link]
- Automatic Differentiation and Neural Networks, [\[Link\]](#)

# 课后作业

- 推荐补充阅读材料