



深度学习框架发展 III

嘉宾：王磊



01 相关背景

- 软件框架
- 机器学习框架
- 深度学习训练框架与推理框架

02 发展历程

- 石器时代
- 青铜时代
- 铁器时代
- 罗马时代
- 工业时代

03 当前发展



- 基于编译器的算子优化
- 分布式训练支持
- 全场景协同
- 训练推理一体化
- 动静态图融合
- 编程接口

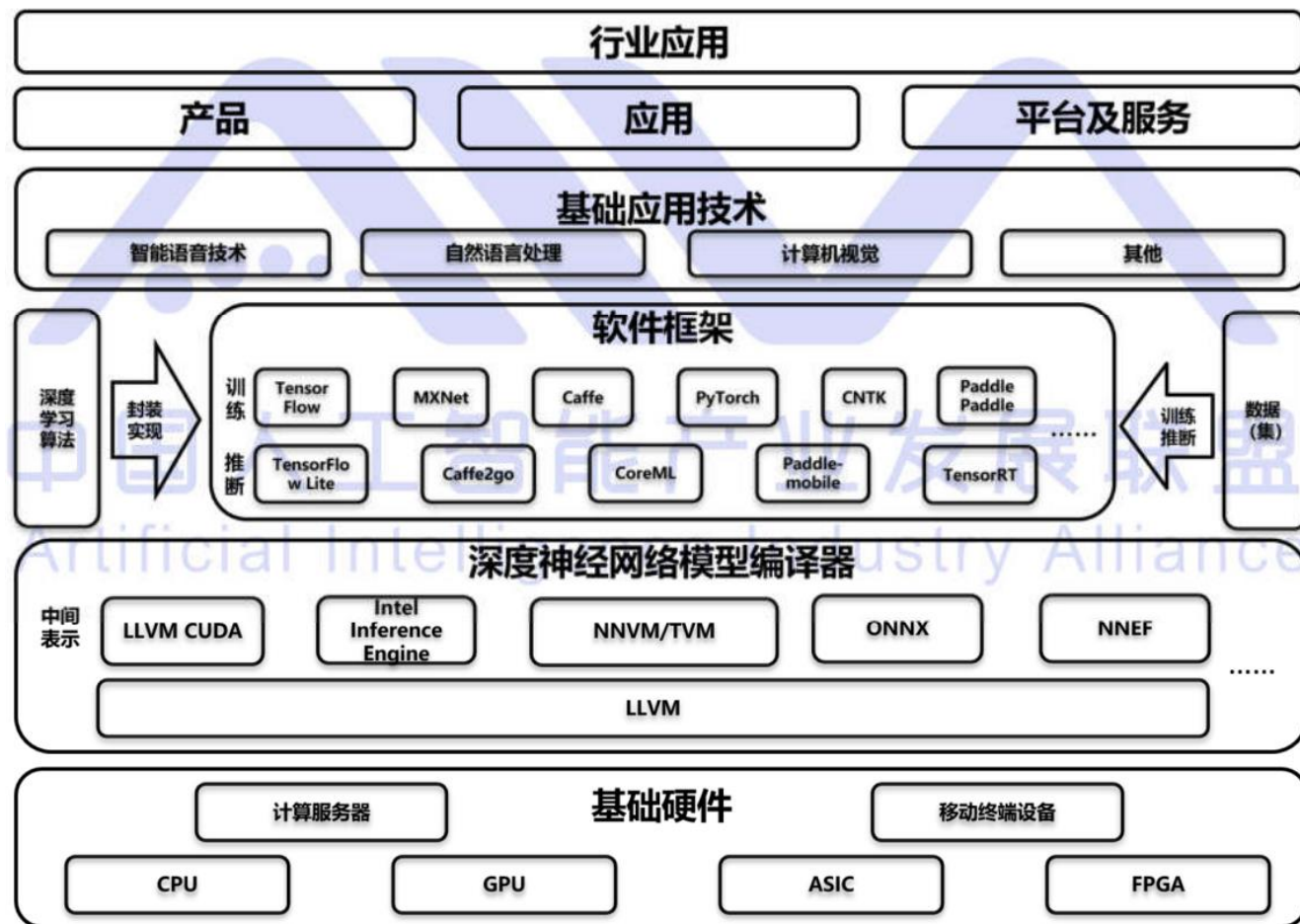


背景——软件框架

3

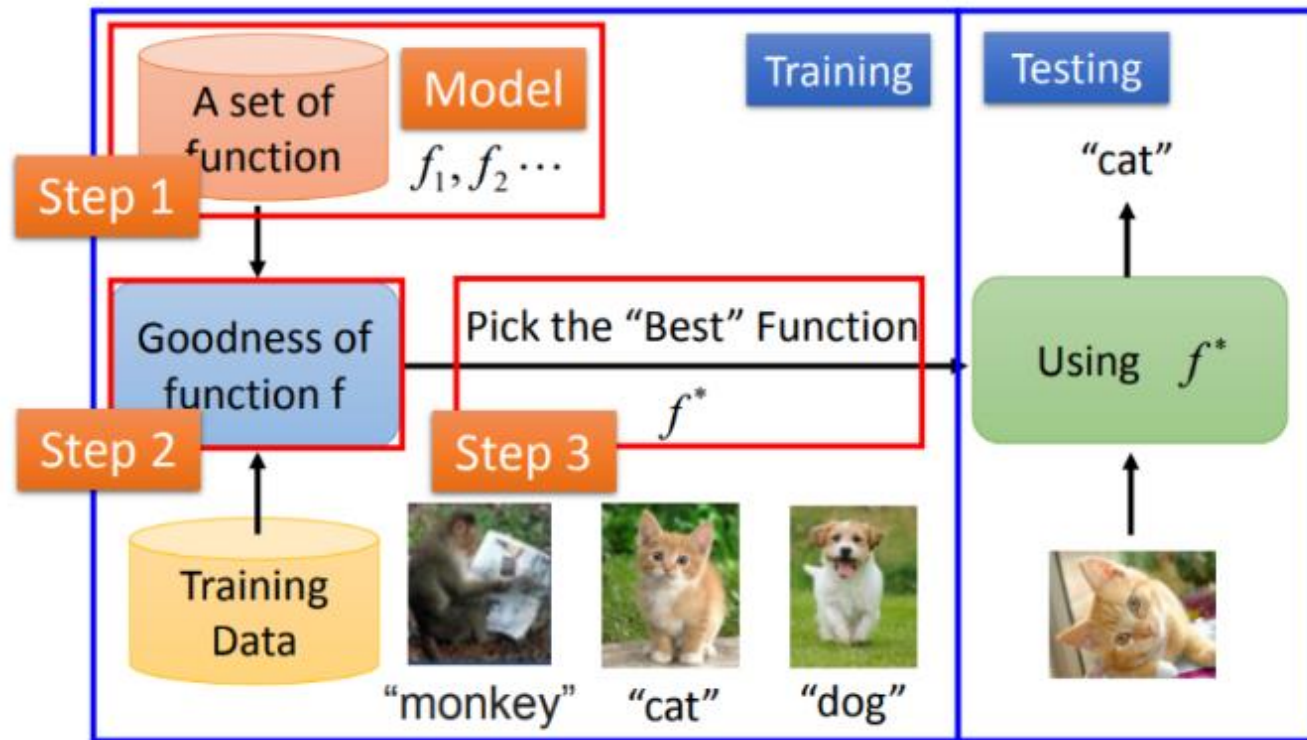
软件框架 (Software Framework)

通常指为了实现某个业界标准或者完成特定基本任务的软件组件规范，也指为了实现某个软件规范时，提供规范所要求的基础功能的软件产品。



机器学习框架 (ML Framework)

机器学习是根据提供的资料，通过使用多种算法寻找function，来自动化各种业务操作并使机器模仿人的行为，但这些算法仅能被ML领域专家所理解，因此，机器学习框架简化了这些复杂的算法，以便可以轻松地在企业应用程序中实现ML。



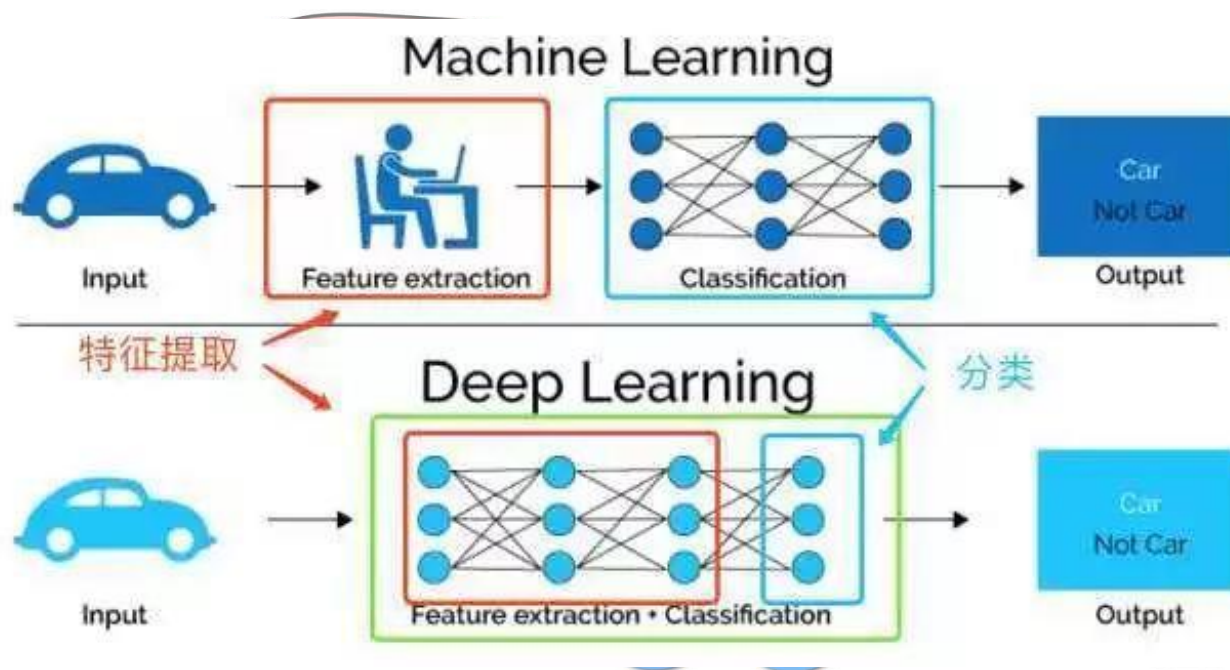
机器学习框架 (ML Framework)

代表性框架Scikit-learn、Apache Mahout、SystemDS、Spark Mllib等

机器学习框架 VS 深度学习框架

机器学习框架包含了深度学习框架

- 对于数据处理的方式不同
- 模型封装的抽象化程度不同
- 针对的群体和项目不同



背景——深度学习训练框架与推理框架



先进编译实验室
Advanced Compiler

6

训练、推理、部署

云端训练框架：主要完成面向海量数据的模型训练任务，对算力要求最高，实际应用中需要采用分布式计算等技术，同时对于工业级模型及稳定性也有特殊要求。

云端推理框架：完成训练模型的优化、云端部署及推断计算等工作，对于效率及并发性等具有特殊要求

端侧推理框架：主要完成训练模型在终端的部署及计算，由于终端功耗、功能、芯片等众多限制，终端推断框架的性能、能耗及自身优化需满足多种限制要求



先进编译实验室
Advanced Compiler



发展历程



先进编译实验室
Advanced Compiler

7



```

1 require 'nn';
2 net = nn.Sequential()
3 net:add(nn.SpatialConvolution(1, 6, 5, 5)) -- 1个图像输
4 net:add(nn.SpatialMaxPooling(2,2,2,2)) -- 1个2x2池
5 net:add(nn.SpatialConvolution(6, 16, 5, 5))
6 net:add(nn.SpatialMaxPooling(2,2,2,2))
7 net:add(nn.View(16*5*5)) -- 将16x5x
8 net:add(nn.Linear(16*5*5, 120)) -- 全连接层
9 net:add(nn.Linear(120, 84))
10 net:add(nn.Linear(84, 10)) -- 10是网络
11
12 print(net:__tostring());
13 mlp=nn.Parallel(2,1); -- iterate over dimension 2
14 mlp:add(nn.Linear(10,3)); -- apply to first slice
15 mlp:add(nn.Linear(10,2)) -- apply to first second sl
16 x = torch.randn(10,2)
17 print(x)
18 print(mlp:forward(x))
19 criterion = nn.ClassNLLCriterion() -- 多分类的负对数似
20 loss = criterion:forward(output, 3) --
21 gradients = criterion:backward(output, 3)
22 gradInput = net:backward(input, gradients)
23 parameters, gradParameters = net:getParameters()
24

```

```

34 d(sub2ind(size(d), D(k), 1)) = 1;
25 -- 定义计算损失和dLoss/dx的闭包。
26 feval = function(x)
27     -- 重置梯度
28     gradParameters:zero()
29     -- 1. 计算每个数据点的输出（对数概率）
30     local output = net:forward(input)
31     -- 2. 根据真实标签，计算这些输出的损失
32     local loss = criterion:forward(output, 3)
33     -- 3. 计算损失对模型输出的导数
34     local dloss_doutput = criterion:backward(output, 3)
35     -- 4. 使用梯度更新权重
36     net:backward(input, dloss_doutput)
37
38     -- 返回 损失 和 相对于权重的损失梯度
39     -- Loss, (gradient of Loss with respect to the weights)
40     return loss, gradParameters
41 end
42 -- 定义 SGD 参数.
43 sgd_params = {
44     learningRate = 1e-2,
45     learningRateDecay = 1e-4,
46     weightDecay = 0,
47     momentum = 0
48 }
49 -- 训练多个批次数据
50 epochs = 1e2
51 losses = {}
52 for i = 1,epochs do
53     -- SGD优化的一个步骤（梯度下降法）
54     _,local_loss = optim.sgd(feval, parameters, sgd_params)
55     -- 累计误差
56     losses[#losses + 1] = local_loss[1]
57 end
58 print(losses[1])
59 print(losses[#losses])
60 print(torch.exp(net:forward(input)))

```


1	layer {	<i># The train/test net protocol buffer definition</i>
2	name: "mnist"	net: "examples/mnist/lenet_train_test.prototxt"
3	type: "Data"	<i># test_iter specifies how many forward passes the test should carry out.</i>
4	transform_param {	<i># In the case of MNIST, we have test batch size 100 and 100 test iterations,</i>
5	scale: 0.00390625	<i># covering the full 10,000 testing images.</i>
6	}	test_iter: 100
7	data_param {	<i># Carry out testing every 500 training iterations.</i>
8	source: "mnist_train_lmdb"	test_interval: 500
9	backend: LMDB	<i># The base learning rate, momentum and the weight decay of the network.</i>
10	batch_size: 64	base_lr: 0.01
11	}	momentum: 0.9
12	top: "data"	weight_decay: 0.0005
13	top: "label"	<i># The Learning rate policy</i>
14	}	lr_policy: "inv"
15	layer {	gamma: 0.0001
16	name: "pool1"	power: 0.75
17	type: "Pooling"	<i># Display every 100 iterations</i>
18	pooling_param {	display: 100
19	kernel_size: 2	<i># The maximum number of iterations</i>
20	stride: 2	max_iter: 10000
21	pool: MAX	<i># snapshot intermediate results</i>
22	}	snapshot: 5000
23	bottom: "conv1"	snapshot_prefix: "examples/mnist/lenet"
24	top: "pool1"	<i># solver mode: CPU or GPU</i>
25	}	solver_mode: GPU
26	

发展历程——铁器时代



先进编译实验室
Advanced Compiler

10

2015~2018年：效率优化、分布式支持、蓬勃发展



PYTORCH



CNTK



先进编译实验室
Advanced Compiler



发展历程——罗马时代



先进编译实验室
Advanced Compiler

11

2018~2019：大模型训练、并行计算、可用性提升



先进编译实验室
Advanced Compiler



当前发展——工业时代



先进编译实验室
Advanced Compiler

12

2020+: 多场景多任务支持、编译层优化、丰富的套件支持

◆ 基于编译器的算子优化

◆ AI+Science

◆ 编程接口优化

◆ 性能领先的模型库

◆ 自动分布式训练

◆ 全场景、异构设备支持

◆ 动态图静态图融合

◆ 训练推理一体化

◆ 小样本训练支持



Advanced Compiler

AI框架发展白皮书（2022年）

<https://syncedreview.com/2020/12/14/a-brief-history-of-deep-learning-frameworks/>

