

复习

- 正则化是深度学习中的一种常用技术，请简述正则化的意义，并列举几种常用的正则化方法；
- 比较Batch Normalization（批量归一化）和Layer Normalization（层归一化）的差异。

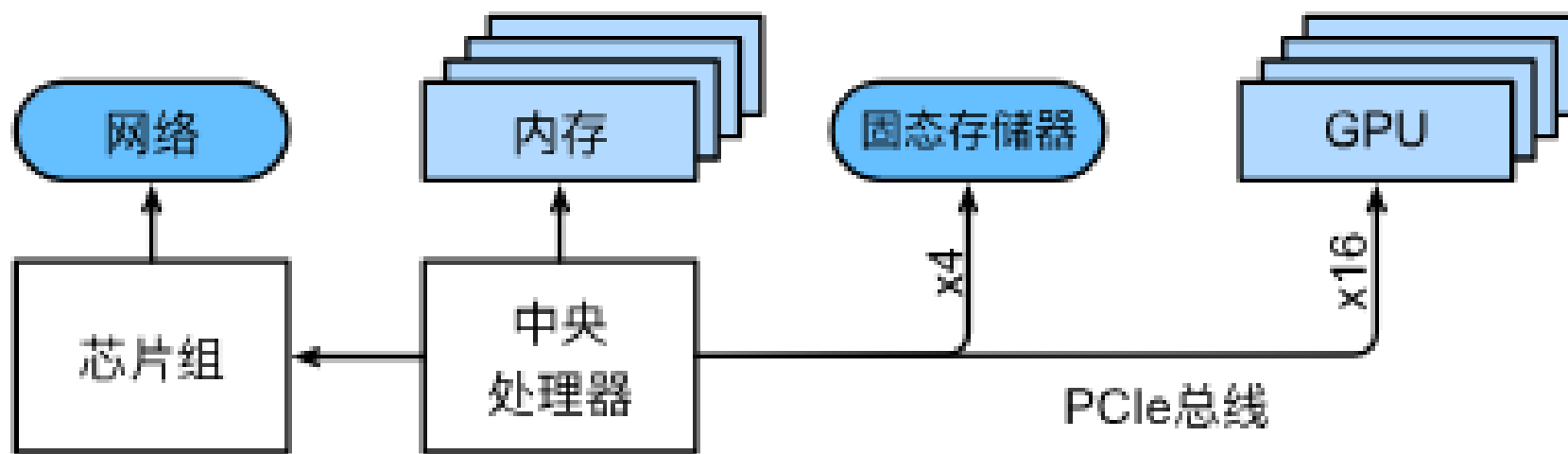


深度学习计算系统

深度学习计算系统

- 深度学习硬件系统
- 深度学习框架
- 深度学习编译
- 深度学习部署
- MLOps

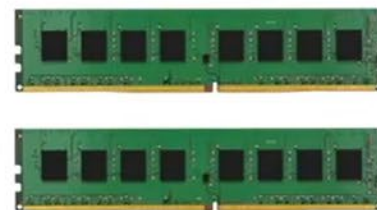
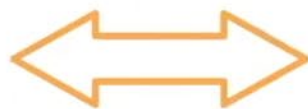
深度学习硬件系统



深度学习硬件系统

Intel i7

0.15 TFLOPS



DDR4

32 GB

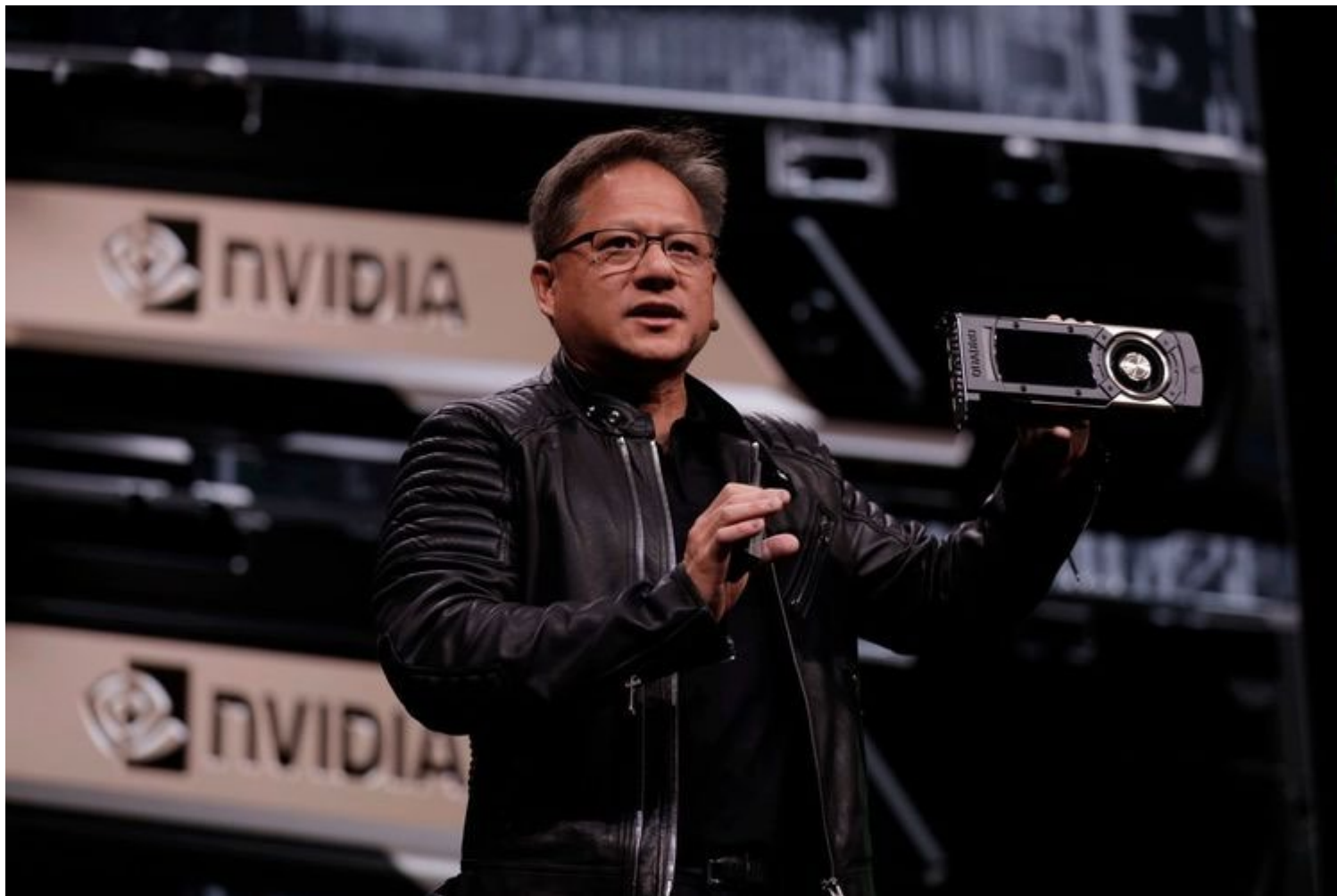


Nvidia Titan X

12 TFLOPS

16 GB

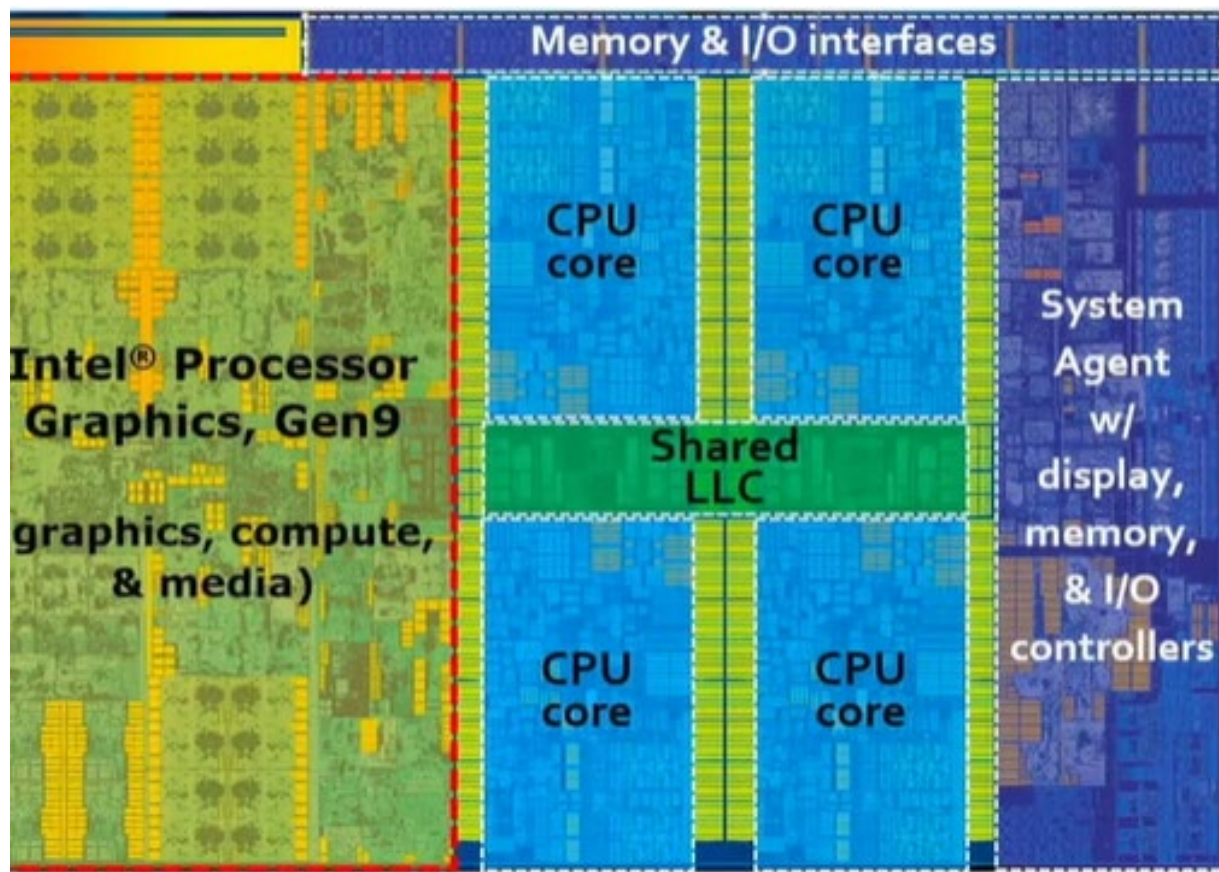
深度学习硬件系统



深度学习硬件系统



深度学习硬件系统

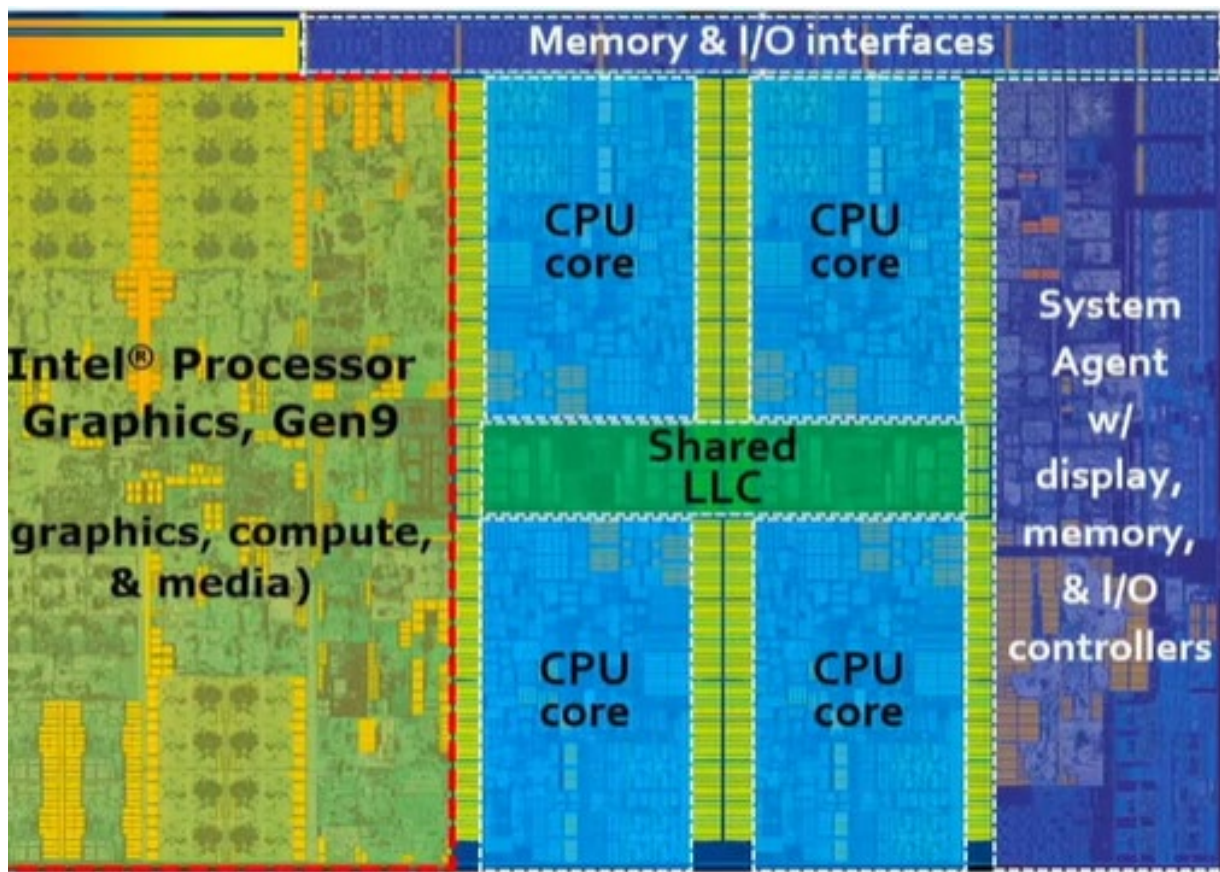


Intel i7-6700K

CPU

- CPU内核
- CPU高速缓存
- 显示处理器（核显）
- 显示、内存、I/O控制器
- 内存、I/O接口

深度学习硬件系统

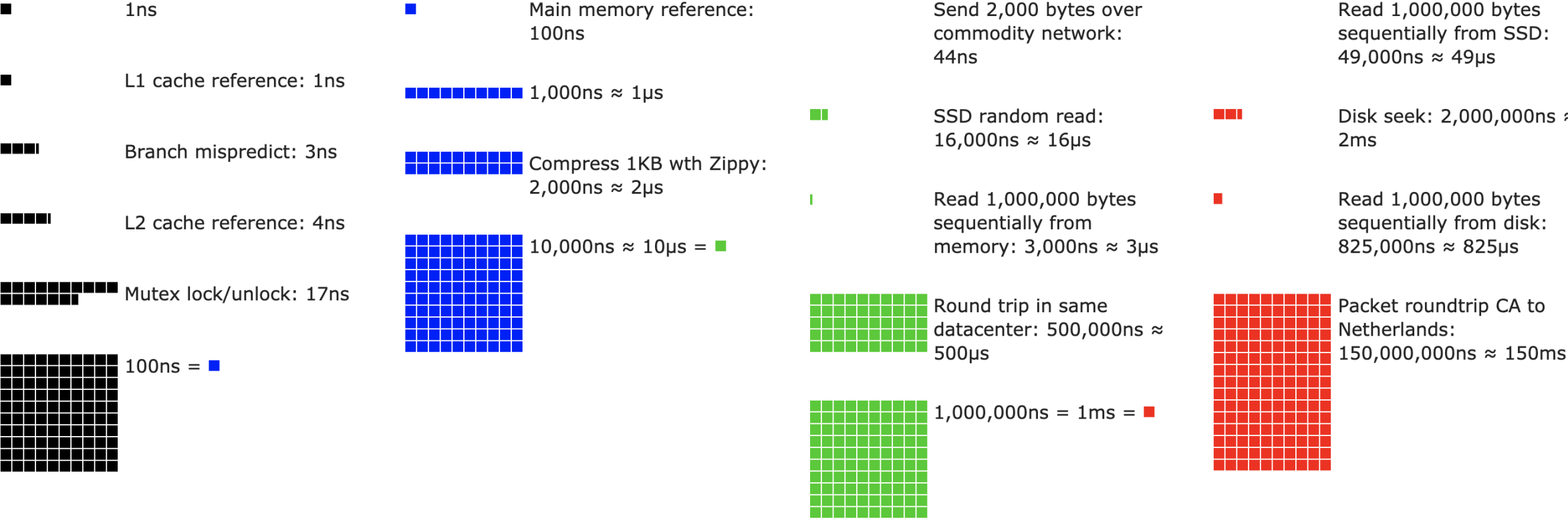


Intel i7-6700K

如何提升CPU利用率？

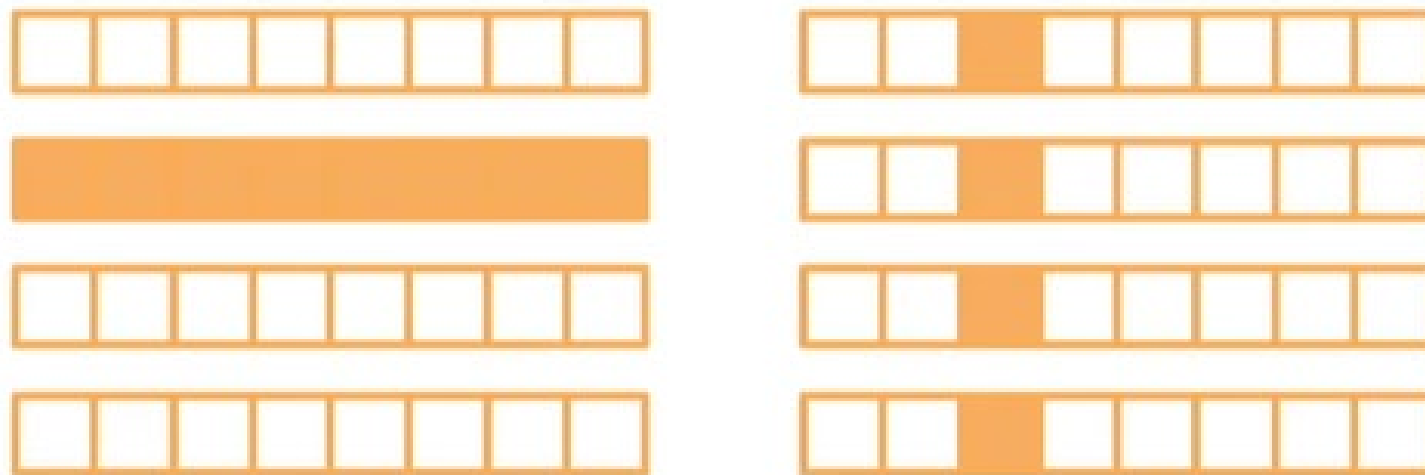
- 在计算 $a+b$ 之前，需要准备数据：
 - 主内存→L3→L2→L1→寄存器；
 - L1访问延时：0.5ns；
 - L2访问延时：7ns；
 - 主内存访问延时：100ns；
- 提升空间和时间的内存本地性：
 - 时间：重用数据，保持数据在缓存；
 - 空间：按序读写数据使得可以预读取。

深度学习硬件系统

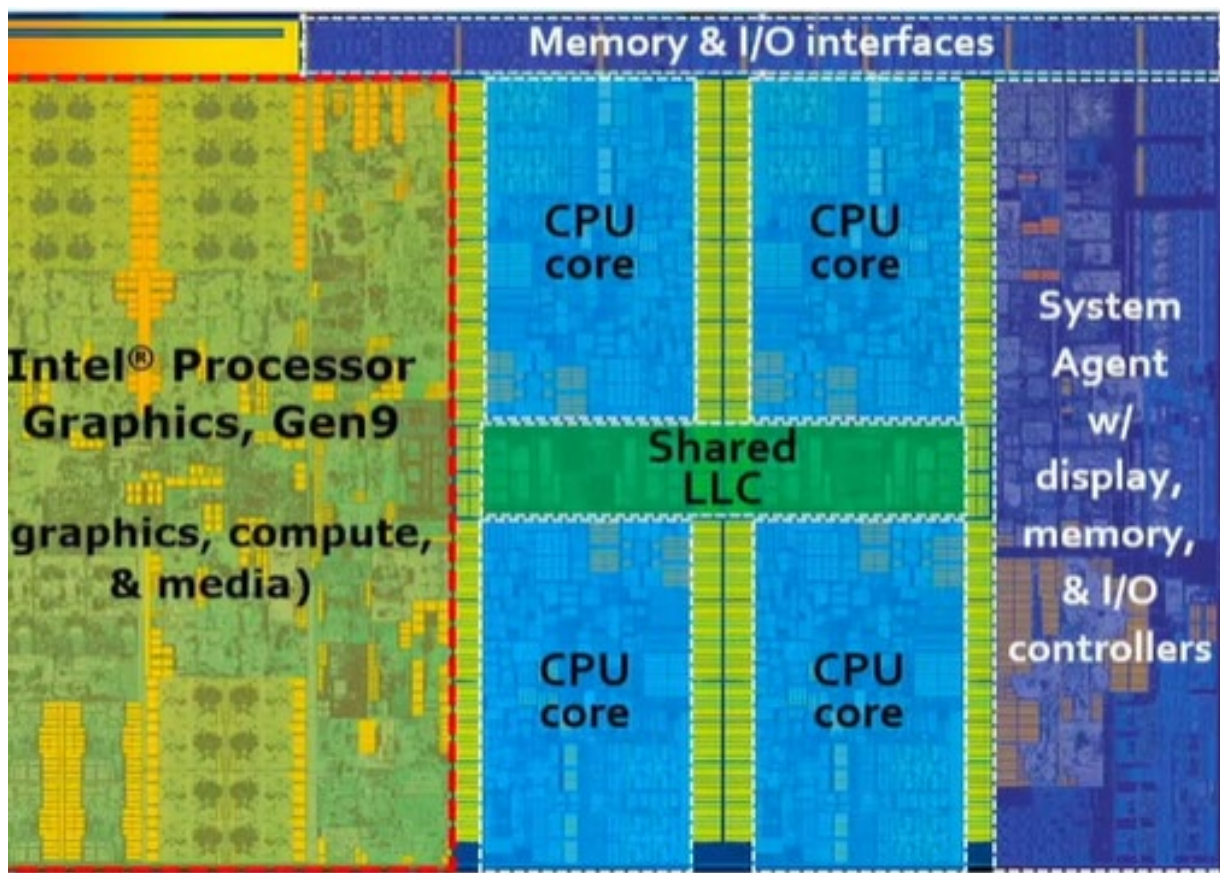


深度学习硬件系统

- 假设矩阵按行存储，那么访问一行比访问一列更快；
 - CPU一次读取64字节（Cache Line）
 - CPU会“聪明地”提前读取下一个（Cache Line）
 - 如果需要按列取矩阵元素，需要进行几次读取？



深度学习硬件系统



Intel i7-6700K

如何提升CPU利用率？

- 另一种策略——并行：
 - 高端CPU通常有几十个核心；
 - 并行利用所有核心；
 - 超线程不一定能提升计算性能，因为它们会共享寄存器，在深度学习计算中，关闭超线程可能还会略微提升性能。

深度学习硬件系统

- 左边比右边慢 (Python) :

```
for i in range(len(a)):
    c[i] = a[i] + b[i]
```

```
c = a + b
```

- 左边调用n次函数，每次调用都有开销；
- 右边容易并行（类似下面的C++实现）：

```
#pragma omp for
for (i=0; i<a.size(); i++) {
    c[i] = a[i] + b[i]
}
```


深度学习硬件系统

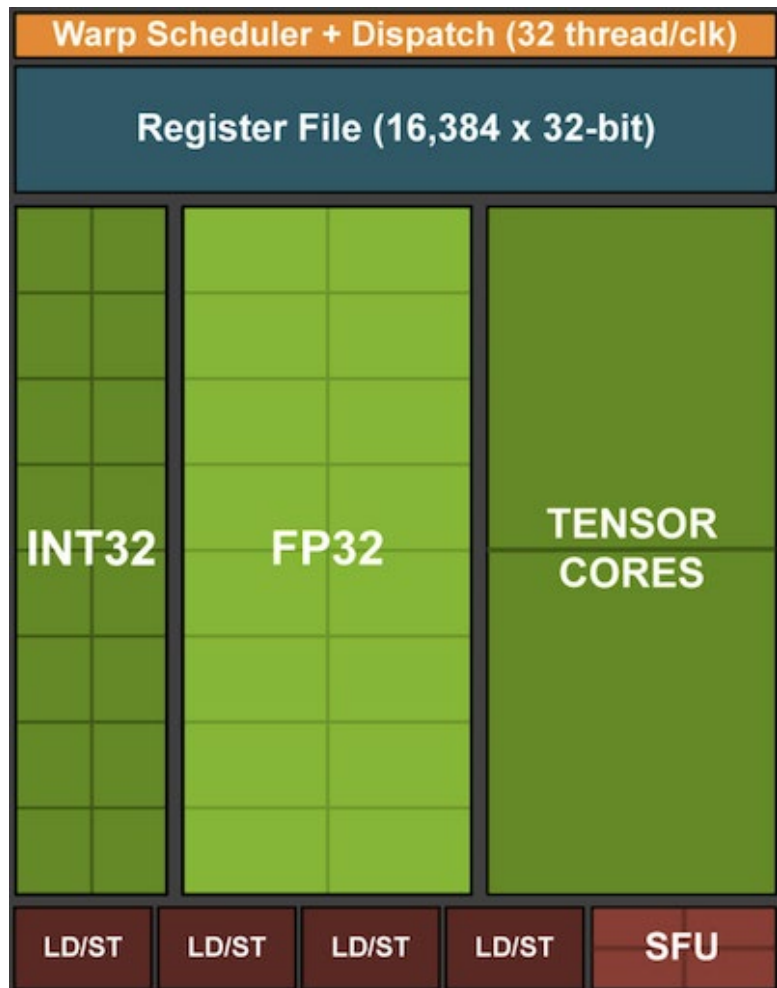


NVIDIA Titan X (Pascal)

GPU

- 流式多处理器 (Stream Multi-processor, SM);
- 每个流式处理单元只能执行简单计算指令（例如浮点运算），但是数量多，并行能力比CPU强。

深度学习硬件系统



GPU

- 新架构的NVIDIA GPU专门设计了张量核（Tensor Core）。

深度学习硬件系统



# 核	6 / 64	2K / 4K
TFLOPS	0.2 / 1	10 / 100
内存大小	32 GB / 1 TB	16 GB / 32 GB
内存带宽	30 GB/s / 100 GB/s	400 GB/s / 1 TB/s
控制流	强	弱

CPU——通用， GPU——特化

深度学习硬件系统

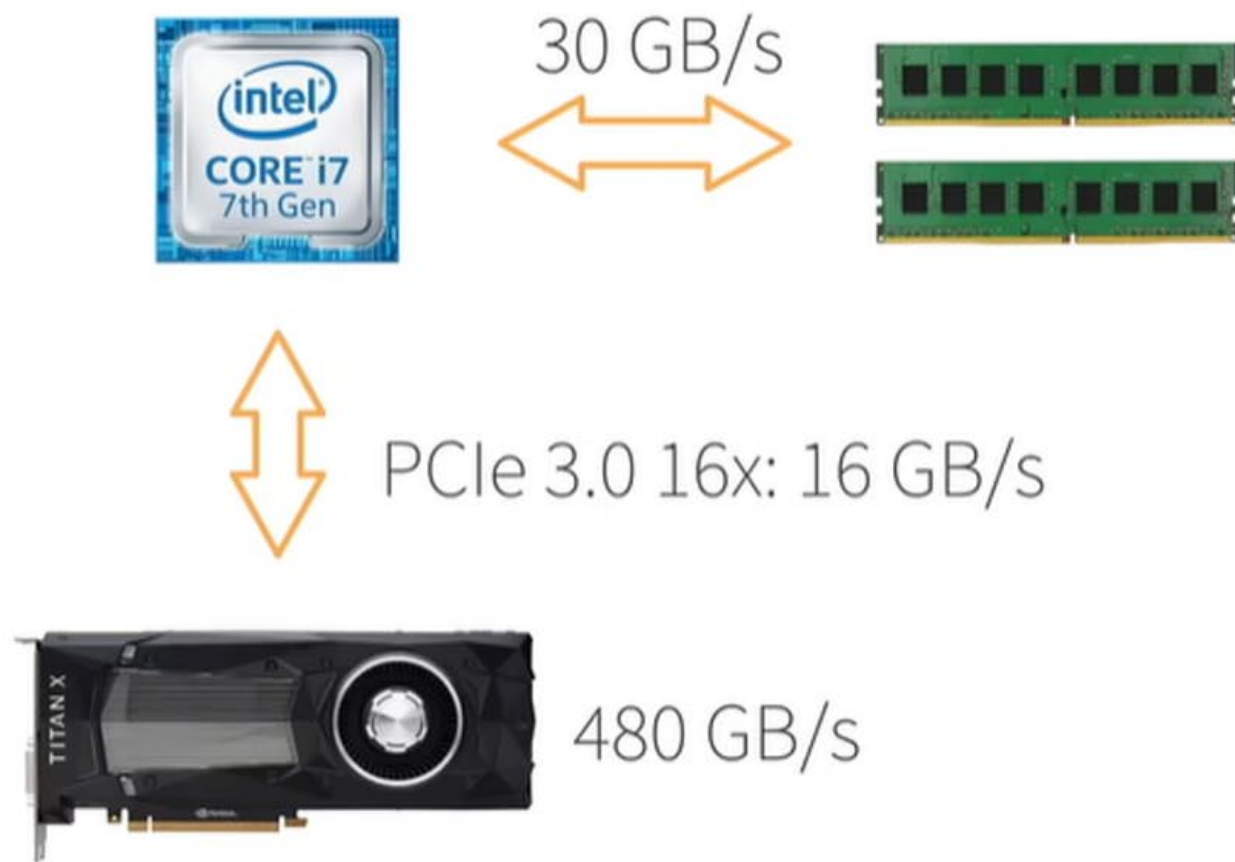


NVIDIA Titan X (Pascal)

提高GPU利用率

- 并行：
使用数千个计算线程（在较大的模型才能很好地体现GPU并行优势）；
- 内存本地性：
缓存小，架构更加简单；
- 少用控制语句：
支持有限，同步开销大。

深度学习硬件系统



- 不要频繁在CPU和GPU之间传数据：带宽限制，同步开销

深度学习硬件系统

更多CPU、GPU

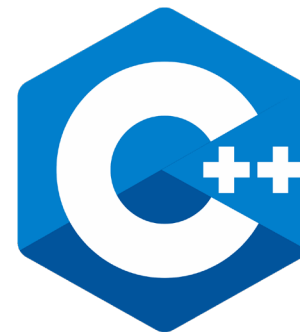
- CPU: AMD, ARM
- GPU: AMD, Intel, ARM, Qualcomm...



CPU/GPU高性能计算编程

- **CPU:**

- C++或其他任何高性能语言
- 编译器成熟



- **GPU:**

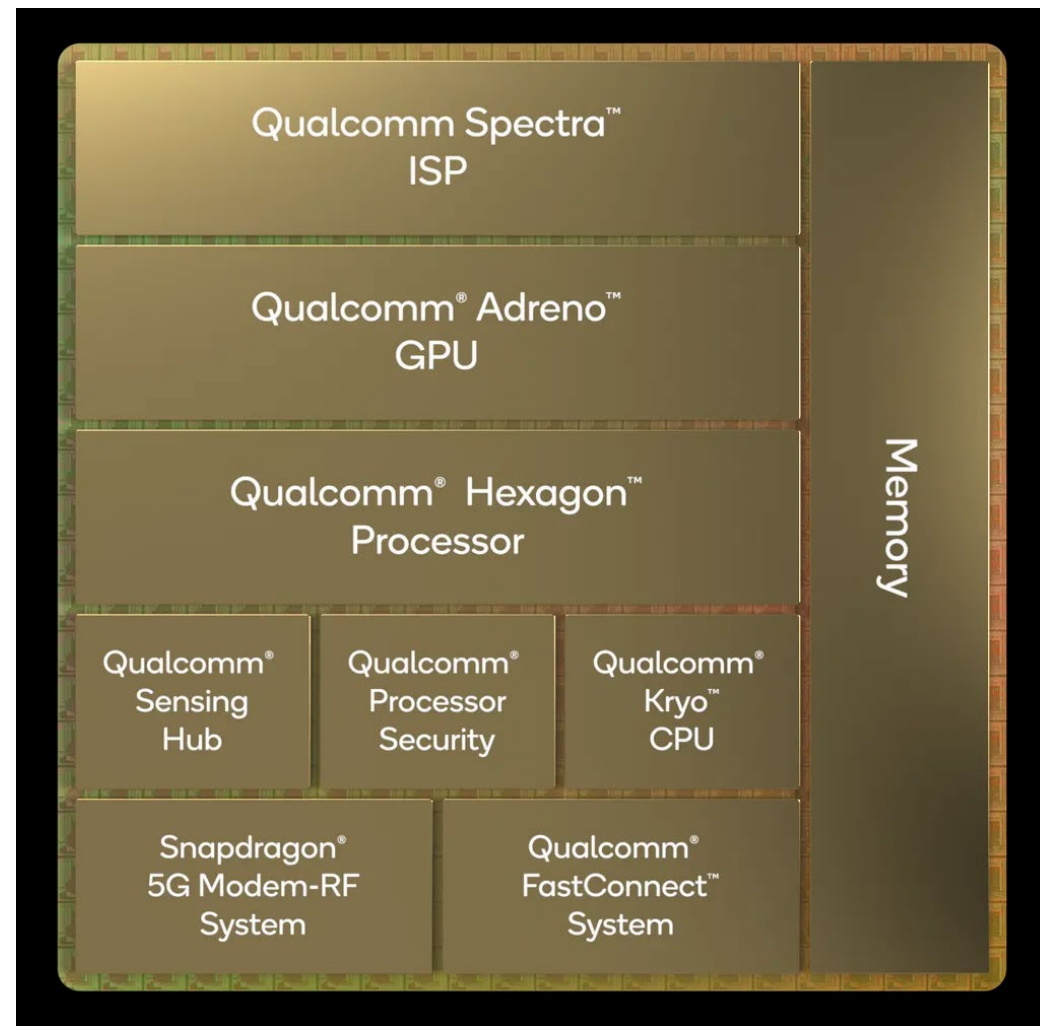
- NVIDIA上用CUDA, 编译器和驱动成熟
- 其他用OpenCL, 质量取决于硬件厂商



深度学习硬件系统

DSP：数字信号处理

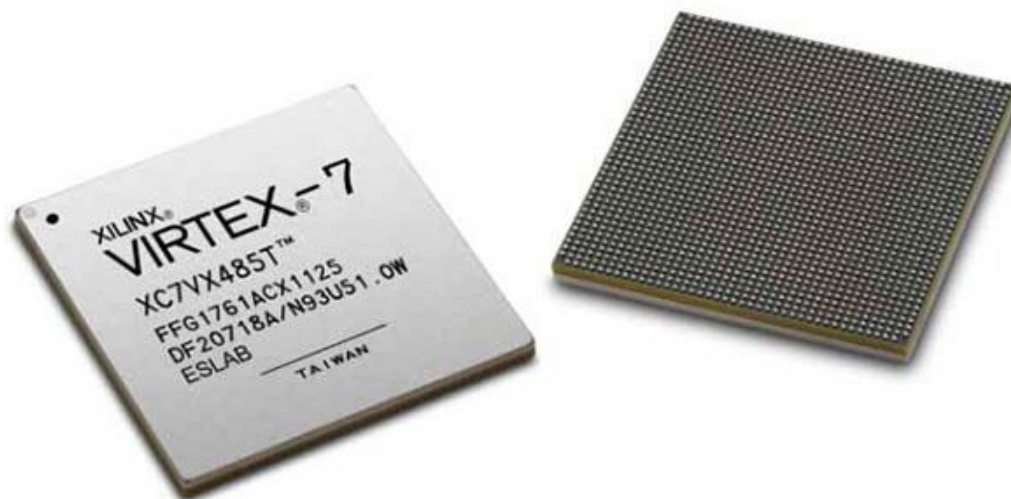
- 为数字信号处理算法设计（滤波、降噪）；
- 低功耗、高性能；
- VLIW: Very long instruction word, 一条指令计算上百次累加；
- 编程和调试困难；
- 编译器质量良莠不齐。



深度学习硬件系统

FPGA：可编程阵列

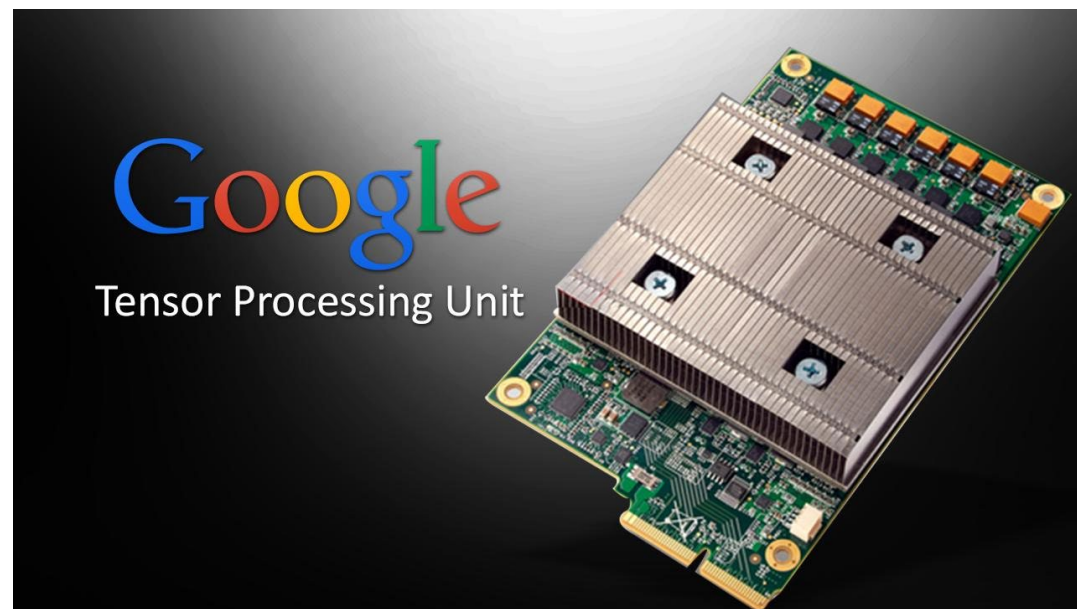
- 有大量可以编程的逻辑单元和可配置连接；
- 可以配置成计算复杂函数：
- 编程语言：VHDL、Verilog；
- 通常比CPU效率更高；
- 工具链质量良莠不齐；
- 一次“编译”需要数小时。



深度学习硬件系统

ASIC：专用集成电路

- 深度学习热门领域：
 - 大公司都在设计制造自己的芯片 (Intel、Qualcomm、Google、Amazon、Facebook.....) ；
- Google TPU是标志性芯片：
 - 能够媲美NVIDIA GPU性能；
 - 在Google大量部署；
 - 核心是systolic array；
- 华为昇腾、寒武纪。



深度学习硬件系统



深度学习硬件系统

征程[®]6 旗舰
Journey

560 TOPS*

专为新一代城区高阶智能驾驶而生
的车载计算方案

城区高阶智驾规模化落地和体验升级，面临复杂嵌入式超级计算机系统技术和极致工程能力的技术挑战，更需要性能高效且资源均衡的计算方案。征程6专为新一代城区高阶智能驾驶推出性能旗舰版本。高集成度的征程6旗舰，算力560TOPS，对BEV Transformer、大规模交互式博弈等先进算法的支持效率领跑业界。征程6旗舰同时提供更均衡的计算资源分配，单芯片可处理感知、决策、规控全流程数据，同时为客户提供灵活的算法优化空间。征程6旗舰的推出将赋能车企全面落地城区高阶智能驾驶，打造差异化用户体验，助力品牌向上。



深度学习硬件系统



并行训练和推理

单机多卡并行

- 一台机器可以安装多个GPU (1-16) ;
- 在训练和推理时, 将一个小批量的计算切分到多个GPU, 来达到加速目的;
- 常用的切分方案:
 - 数据并行;
 - 模型并行;
 - 通道并行 (数据+模型) 。

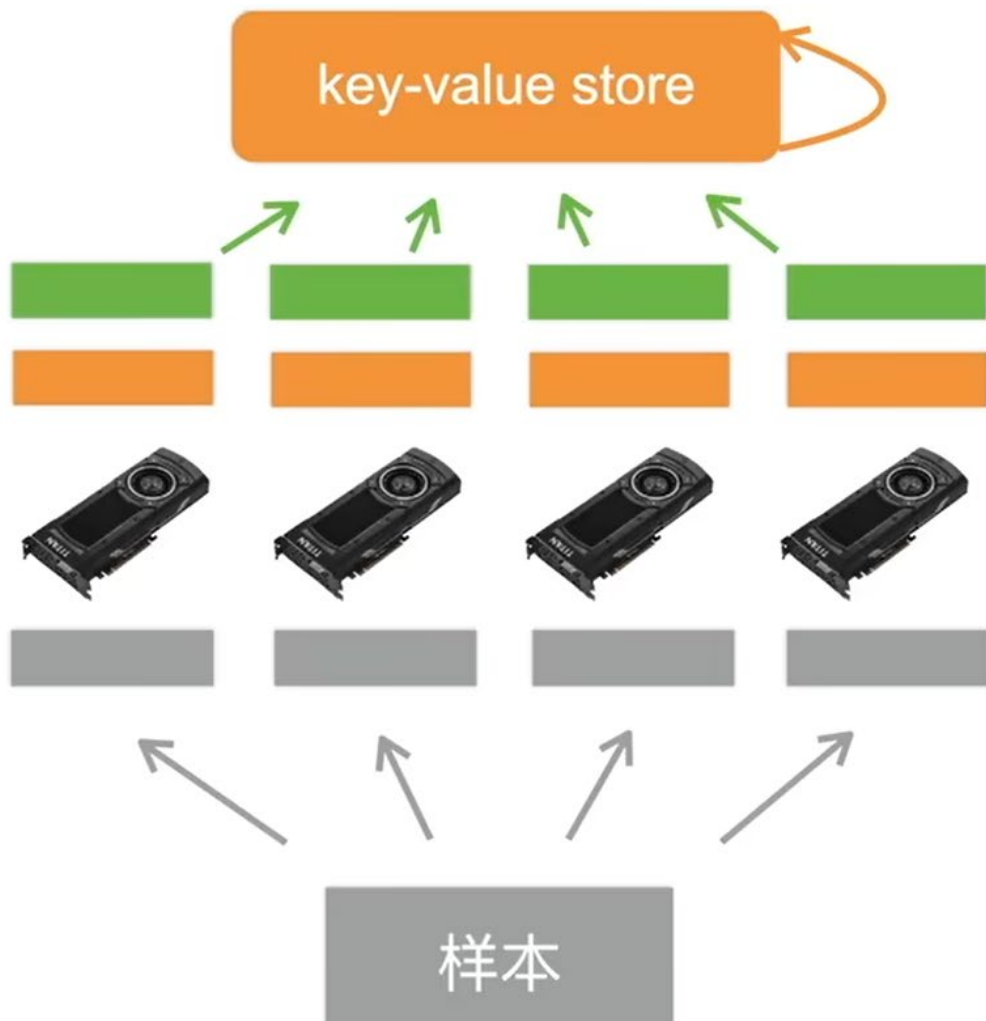
并行训练和推理

数据并行 VS 模型并行

- 数据并行：
将小批量分成 n 块，每个GPU拿到完整的参数计算1块数据的梯度；
 - 通常性能更好；
- 模型并行：
将模型分成 n 块，每个GPU拿到1块模型计算它的前向和反向结果；
 - 通常用于规模很大的模型。

并行训练和推理

数据并行

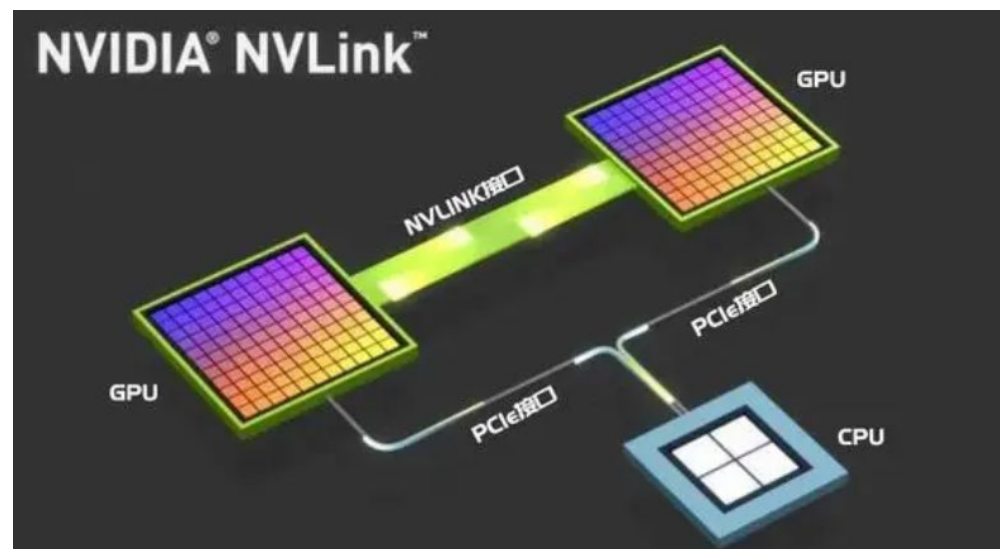


1. 读一个数据块
2. 拿回参数
3. 计算梯度
4. 发出梯度
5. 更新梯度

并行训练和推理

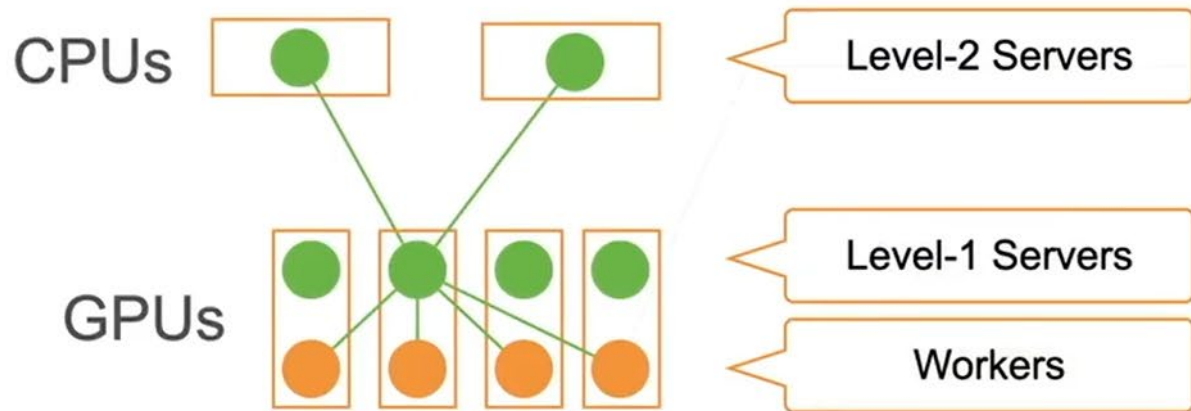
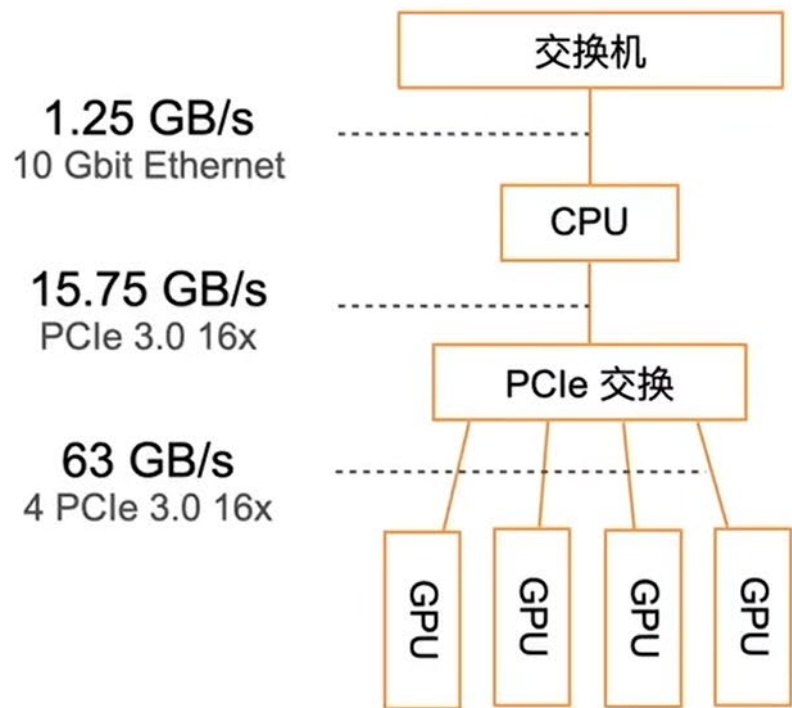
数据并行——PCIe V.S. NVLink

- PCIe：一种基于串行通信的总线标准，采用点对点连接，支持多个设备同时进行高速数据传输；
 - PCIe 4.0的带宽约为32GB/s；
- NVLink：NVidia开发的一种高速连接技术，主要用于连接GPU和其他组件，采用并行通信，较PCIe能提供更高的带宽和更低的延迟；
 - NVLink 2.0的带宽高达300GB/s。



并行训练和推理

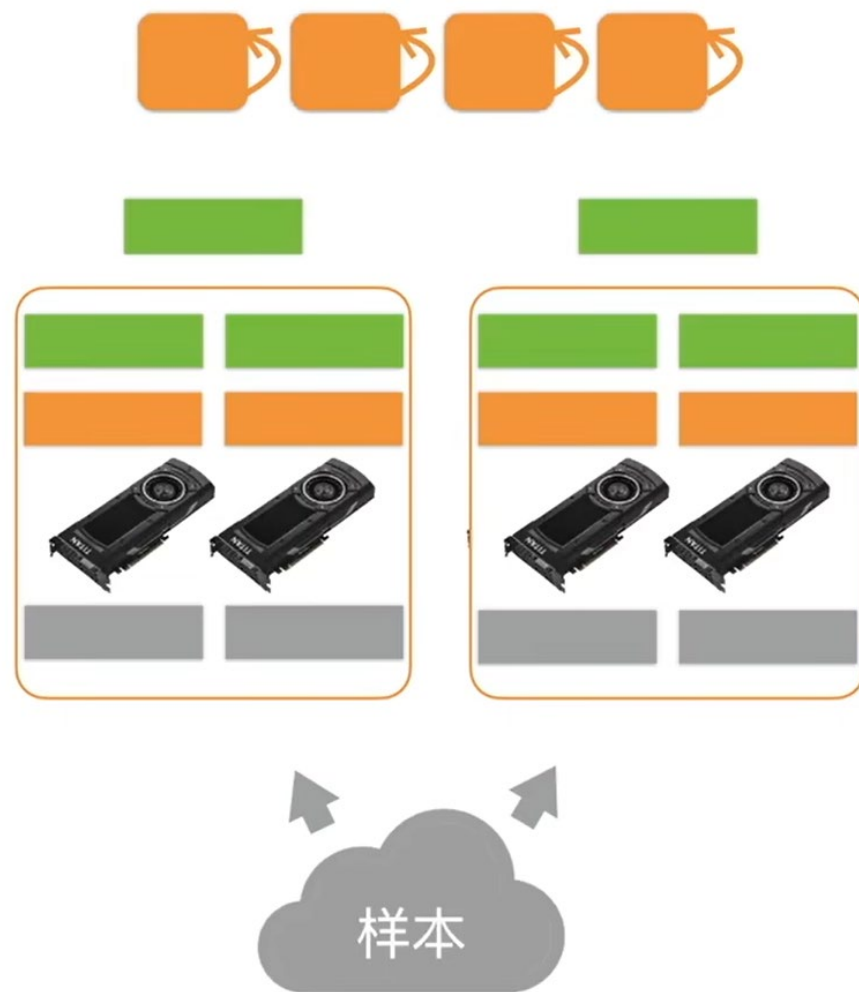
多机多卡并行



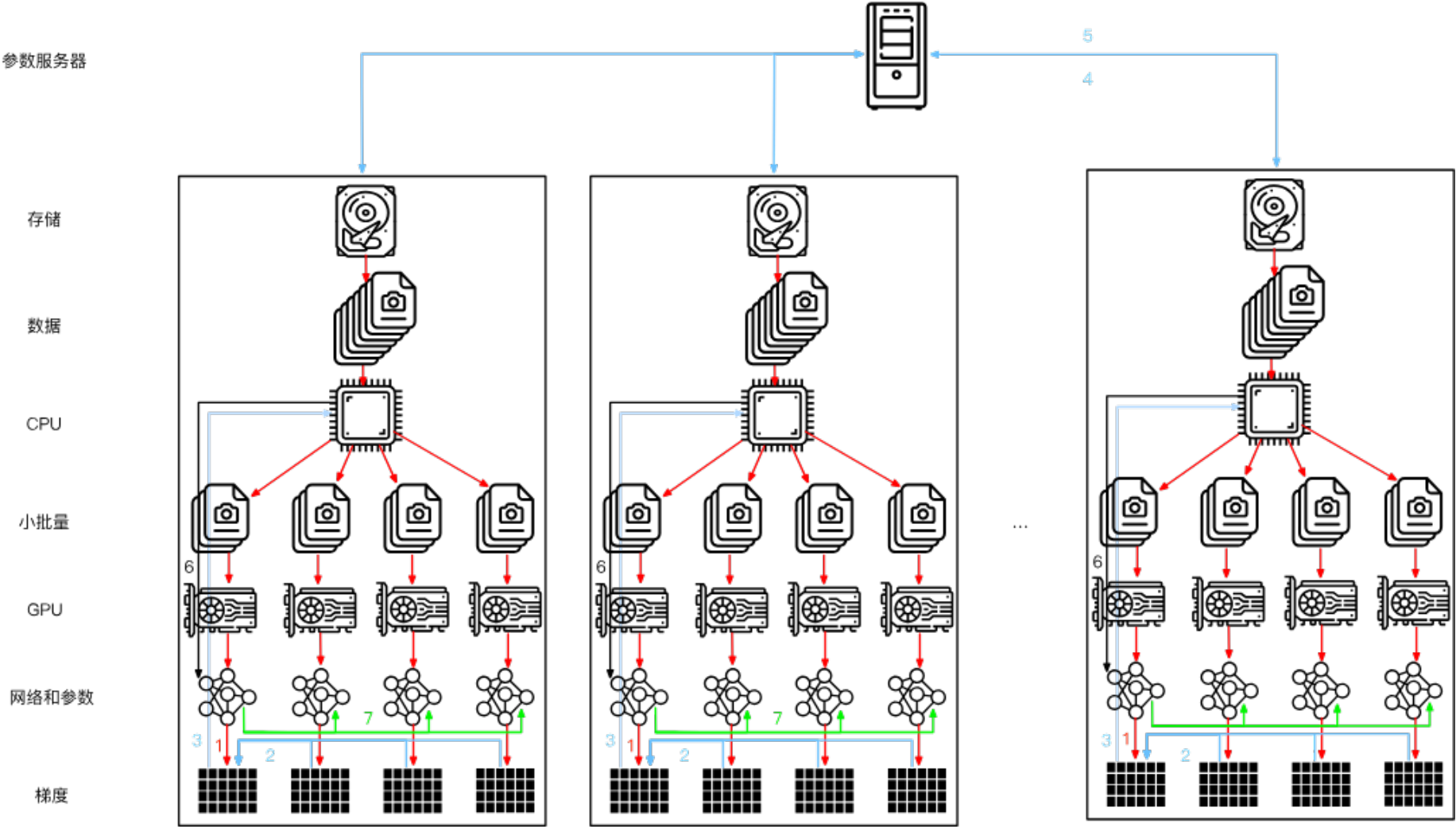
并行训练和推理

多机多卡并行

- 每台计算服务器读取数据的一块；
- 进一步切分数据到每块GPU；
- 每个worker从参数服务器获取模型参数；
- 复制参数到每块GPU；
- 每块GPU计算梯度；
- 将所有GPU上的梯度求和；
- 梯度传回参数服务器；
- 每台服务器更新参数。

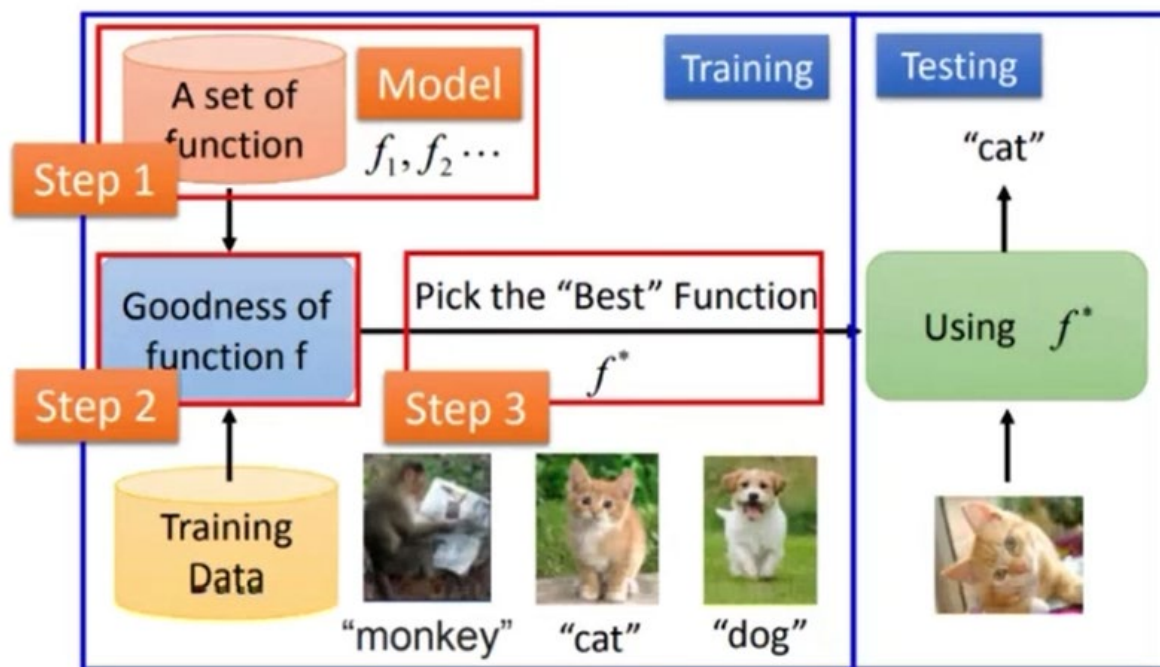


并行训练和推理



深度学习框架

- 软件框架（Software Framework）：为了实现某个业界标准，或完成特定基本任务的软件组件规范；
- 深度学习存在大量重复的基础算法、流程——重复工作。



深度学习框架——发展历程



深度学习框架——发展历程

“石器时代”

~2012年：API不友好，没有GPU支持，尚不完善



MATLAB是由美国MathWorks公司出品的一种用于算法开发、数据分析以及数值计算的高级技术计算语言和交互式环境。



OpenNN 开发于 2003 年在国际工程数值方法中心的名为 RAMFLOOD的项目中是一个使用C++编写的开源类库。



Torch是Facebook的开源机器学习库、科学计算框架和基于Lua编程语言的脚本语言，是PyTorch的直系祖先。

```

1 function [W1, W5, Wo] = MnistConv(W1, W5, Wo, X, D)
2 %
3 alpha = 0.01;
4 beta = 0.95;
5 momentum1 = zeros(size(W1));
6 momentum5 = zeros(size(W5));
7 momentum0 = zeros(size(Wo));
8 N = length(D);
9 bsize = 100;
10 blist = 1:bsize:(N-bsize+1);
11 % One epoch Loop
12 for batch = 1:length(blist)
13     dW1 = zeros(size(W1));
14     dW5 = zeros(size(W5));
15     dWo = zeros(size(Wo));
16     % Mini-batch Loop
17     %
18     begin = blist(batch);
19     for k = begin:begin+bsize-
20         % Forward pass = inference
21         %
22         x = X(:, :, k);           % Input,          28x28
23         y1 = Conv(x, W1);          % Convolution,  20x20x20
24         y2 = ReLU(y1);
25         y3 = Pool(y2);             % Pool,         10x10x20
26         y4 = reshape(y3, [], 1);   %               2000
27         v5 = W5*y4;                % ReLU,         360
28         y5 = ReLU(v5);             %
29         v = Wo*y5;                 % Softmax,      10
30         % Softmax loss

```

```

34     d(sub2ind(size(d), D(k), 1)) = 1;
35     % Backpropagation
36     %
37     e = d - y;                    % Output Layer
38     delta = e;
39     e5 = Wo' * delta;              % Hidden(ReLU) Layer
40     delta5 = (y5 > 0) .* e5;
41     e4 = W5' * delta5;            % Pooling Layer
42     e3 = reshape(e4, size(y3));
43     e2 = zeros(size(y2));
44     W3 = ones(size(y2)) / (2*2);
45     for c = 1:20
46         e2(:, :, c) = kron(e3(:, :, c), ones([2 2])) .* W3(:, :, c);
47     end
48     delta2 = (y2 > 0) .* e2;      % ReLU Layer
49     delta1_x = zeros(size(W1));   % Convolutional Layer
50     for c = 1:20
51         delta1_x(:, :, c) = conv2(x(:, :, c), rot90(delta2(:, :, c), 2), 'valid');
52     end
53     dW1 = dW1 + delta1_x;
54     dW5 = dW5 + delta5*y4';
55     dWo = dWo + delta *y5';
56 end
57 % Update weights
58 %
59 dW1 = dW1 / bsize;
60 dW5 = dW5 / bsize;
61 dWo = dWo / bsize;
62 momentum1 = alpha*dW1 + beta*momentum1;
63 W1 = W1 + momentum1;
64 momentum5 = alpha*dW5 + beta*momentum5;
65 W5 = W5 + momentum5;
66 momentum0 = alpha*dWo + beta*momentum0;
67 Wo = Wo + momentum0;
68 end

```

```

1  require 'nn';
2  net = nn.Sequential()
3  net:add(nn.SpatialConvolution(1, 6, 5, 5)) -- 1个图像输入channel, 6个输出channels, 5x5卷积核
4  net:add(nn.SpatialMaxPooling(2,2,2,2))    -- 1个2x2窗口的最大池化操作
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16x5x5的3D张量变为1D的16*5*5
8  net:add(nn.Linear(16*5*5, 120))           -- 全连接层(输入和权重之间的矩阵乘法)
9  net:add(nn.Linear(120, 84))
10 net:add(nn.Linear(84, 10))                -- 10是网络输出的数量net:add(nn.LogSoftMax())
11
12 print(net:__tostring());
13 mlp=nn.Parallel(2,1);    -- iterate over dimension 2 of input
14 mlp:add(nn.Linear(10,3)); -- apply to first slice
15 mlp:add(nn.Linear(10,2)) -- apply to first second slice
16 x = torch.randn(10,2)
17 print(x)
18 print(mlp:forward(x))
19 criterion = nn.ClassNLLCriterion() -- 多分类的负对数似然准则
20 loss = criterion:forward(output, 3) --
21 gradients = criterion:backward(output, 3)
22 gradInput = net:backward(input, gradients)
23 parameters, gradParameters = net:getParameters()

```

深度学习框架——发展历程

“青铜时代”

2013~2015年：CPU多核、GPU支持、不同编程风格

Caffe

Caffe是以C++/CUDA代码为主的深度学习框架，需要进行编译安装，支持命令行、Python和MATLAB接口，单机多卡、多机多卡等都可以很方便的使用

theano

Theano是一个Python库和优化编译器的开源项目，用于操作和评估数学表达式，尤其是矩阵值表达式。其计算是使用NumPy语法表示，并且编译后可在CPU/GPU架构上高效运行



Chainer是一个开源的深度学习框架，完全在 NumPy 和 CuPy Python 库的基础上用 Python 编写,因其采用边运行边定义 方案以及在大型系统上的性能而著称


```

1 layer {
2   name: "mnist"
3   type: "Data"
4   transform_param {
5     scale: 0.00390625
6   }
7   data_param {
8     source: "mnist_train_lmdb"
9     backend: LMDB
10    batch_size: 64
11  }
12  top: "data"
13  top: "label"
14 }
15 layer {
16   name: "pool1"
17   type: "Pooling"
18   pooling_param {
19     kernel_size: 2
20     stride: 2
21     pool: MAX
22   }
23   bottom: "conv1"
24   top: "pool1"
25 }

```

```

# The train/test net protocol buffer definition
net: "examples/mnist/lenet_train_test.prototxt"
# test_iter specifies how many forward passes the test should carry out.
# In the case of MNIST, we have test batch size 100 and 100 test iterations,
# covering the full 10,000 testing images.
test_iter: 100
# Carry out testing every 500 training iterations.
test_interval: 500
# The base learning rate, momentum and the weight decay of the network.
base_lr: 0.01
momentum: 0.9
weight_decay: 0.0005
# The Learning rate policy
lr_policy: "inv"
gamma: 0.0001
power: 0.75
# Display every 100 iterations
display: 100
# The maximum number of iterations
max_iter: 10000
# snapshot intermediate results
snapshot: 5000
snapshot_prefix: "examples/mnist/lenet"

```

深度学习框架——发展历程

“铁器时代”

2015~2018年：效率优化、分布式支持、蓬勃发展



深度学习框架——发展历程

“罗马时代”

2018~2019年：大规模训练、并行计算、可用性提升



深度学习框架——发展历程

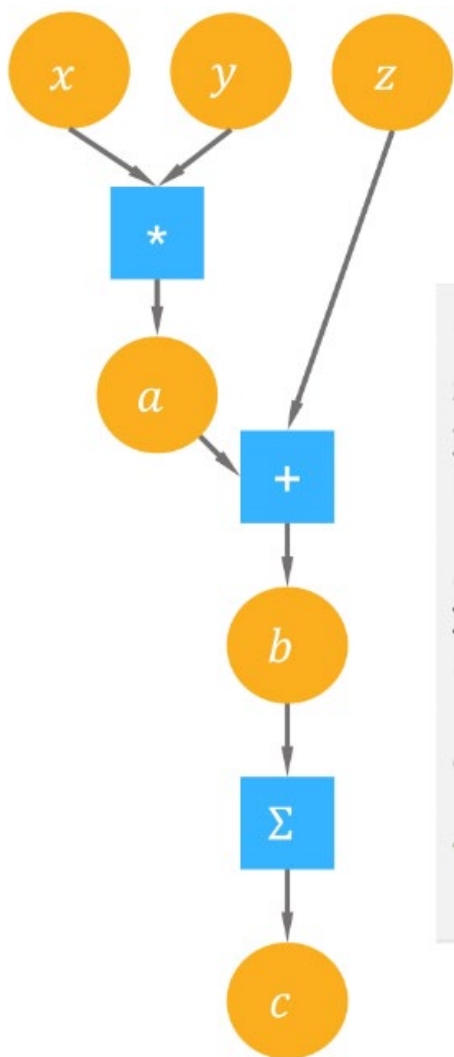
“工业时代”

2020年~：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
- 编程接口优化
- 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



深度学习框架



TensorFlow Program

```
import tensorflow as tf

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)
```



NEW ANNOUNCEMENTS

Catch up on the latest technical insights and tools from the PyTorch community.

[Read More](#) >

MEMBERSHIP AVAILABLE

Join the Membership that fits your goals.

[Join](#)

PYTORCH 2.2


PyTorch 2.2 offers ~2x performance improvement.

[Learn More](#)[PyTorch](#)[torchaudio](#)[torchtext](#)[torchvision](#)[torcharrow](#)[TorchData](#)[TorchRec](#)[TorchServe](#)[PyTorch on XLA Devices](#)

experiences
stability, and

深度学习框架——框架开发

[README](#) [Code of conduct](#) [Apache-2.0 license](#) [Security](#)



TensorFlow

python 3.9 | 3.10 | 3.11 | 3.12

pypi package 2.16.1

DOI 10.5281/zenodo.4724125

openssf best practices passing

openssf scorecard 8.2

oss-fuzz coverage failing

oss-fuzz build failing

OSSRank #9 (Top 1%)

Contributor Covenant v1.4 adopted

TF Official Continuous 8 passed, 0 failed

TF Official Nightly 13 passed, 3 failed

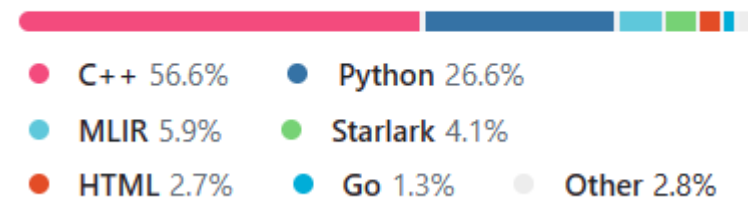
Documentation

api reference

[TensorFlow](#) is an end-to-end open source platform for machine learning. It has a comprehensive, flexible ecosystem of [tools](#), [libraries](#), and [community](#) resources that lets researchers push the state-of-the-art in ML and developers easily build and deploy ML-powered applications.


TensorFlow was originally developed by researchers and engineers working within the Machine Intelligence team at

Languages



深度学习框架——框架开发

[README](#) [Code of conduct](#) [License](#) [Security](#)



PyTorch

PyTorch is a Python package that provides two high-level features:

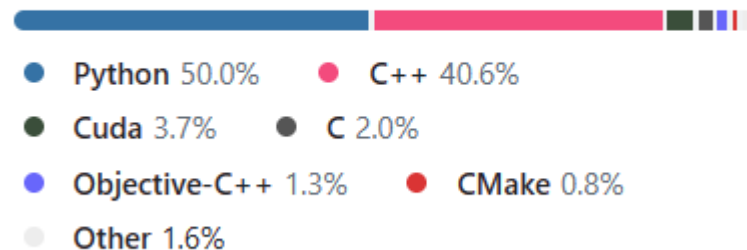
- Tensor computation (like NumPy) with strong GPU acceleration
- Deep neural networks built on a tape-based autograd system

You can reuse your favorite Python packages such as NumPy, SciPy, and Cython to extend PyTorch when needed.

Our trunk health (Continuous Integration signals) can be found at hud.pytorch.org.

- [More About PyTorch](#)
 - [A GPU-Ready Tensor Library](#)
 - [Dynamic Neural Networks: Tape-Based Autograd](#)
 - [Python First](#)
 - [Imperative Experiences](#)

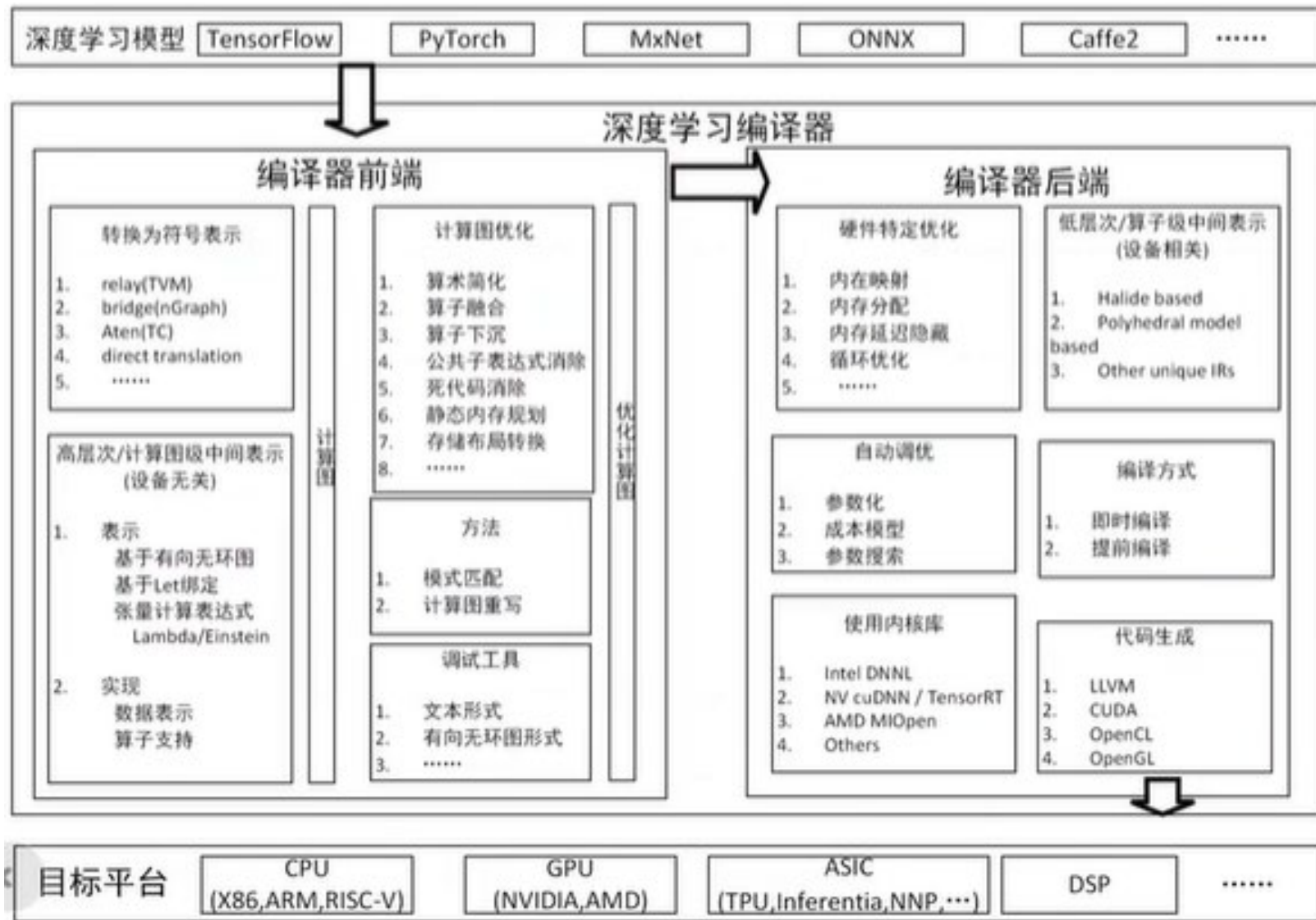
Languages



深度学习框架——框架开发

- 现代C++
- CUDA
- C++代码封装（提供Python端接口）
- CMake——自动构建系统

深度学习编译系统



深度学习部署

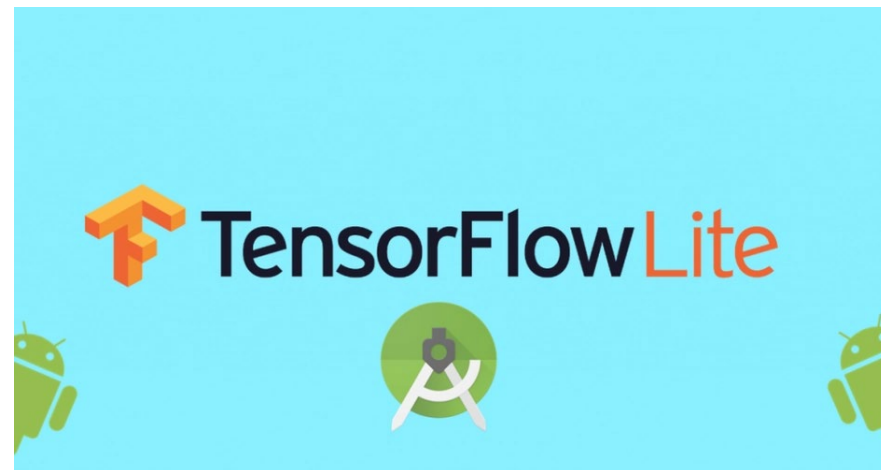
部署方式:

- 训练框架部署：如使用完整的TensorFlow、PyTorch等；
 - 需要安装整个框架；
 - 推理性能差；
 - 很多冗余功能；
 - 内存占用大。

深度学习部署

部署方式：

- 训练框架部署引擎：
TensorFlow Serving、
TensorFlow Lite、TorchServe等；
 - 只支持自己框架训练的模型；
 - 支持的硬件、OS有限。



TORCHSERVE

TorchServe is a performant, flexible and easy to use tool for serving PyTorch

What's going on in TorchServe?

- High performance Llama 2 deployments with AWS Inferentia2 using TorchServe
 - Naver Case Study: Transition From High-Cost GPUs to Intel CPUs and on-prem
 - Run multiple generative AI models on GPU using Amazon SageMaker mlflow
- 75% in inference costs

深度学习部署

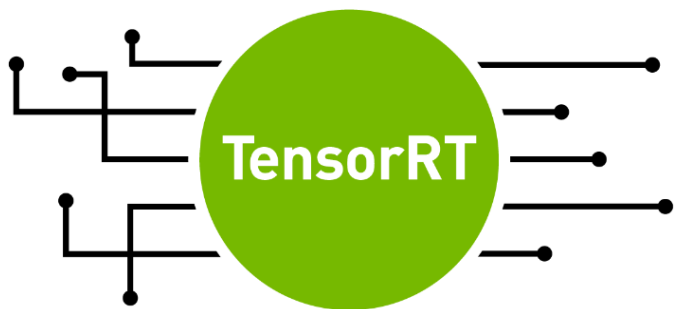
部署方式：

- 手动模型重构：
手写C/C++代码，实现计算图，导入权重；
 - “从头造轮子”；
 - 需要开发者对模型有充分的了解；
 - 有较高的技术难度。

深度学习部署

部署方式：

- 深度学习推理引擎：
ONNX、OpenVINO、NVIDIA Tensor RT，移动端NCNN；
 - 支持加载主流框架的模型文件；
 - 性能良好，能够满足不同环境下的应用需求。

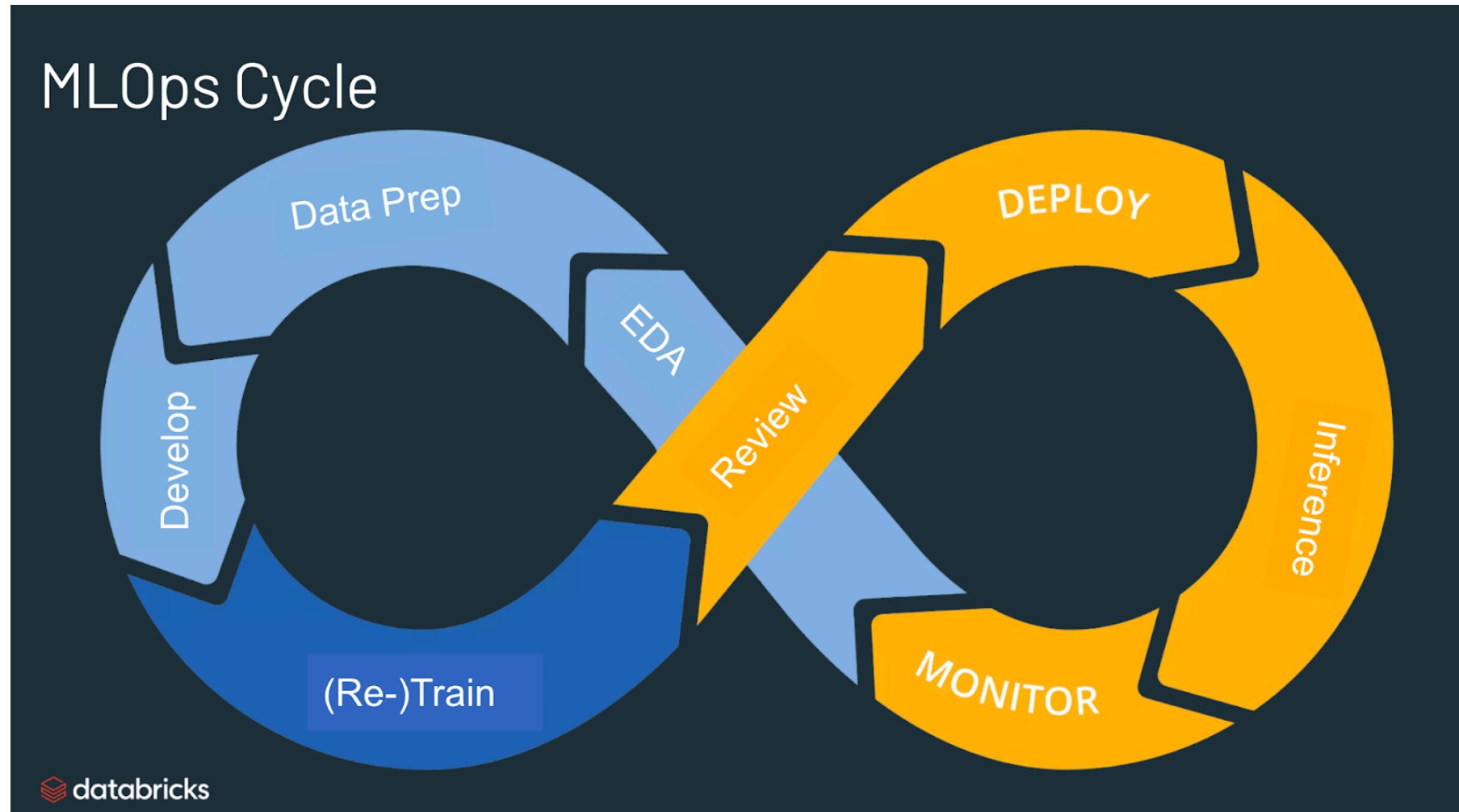


深度学习部署

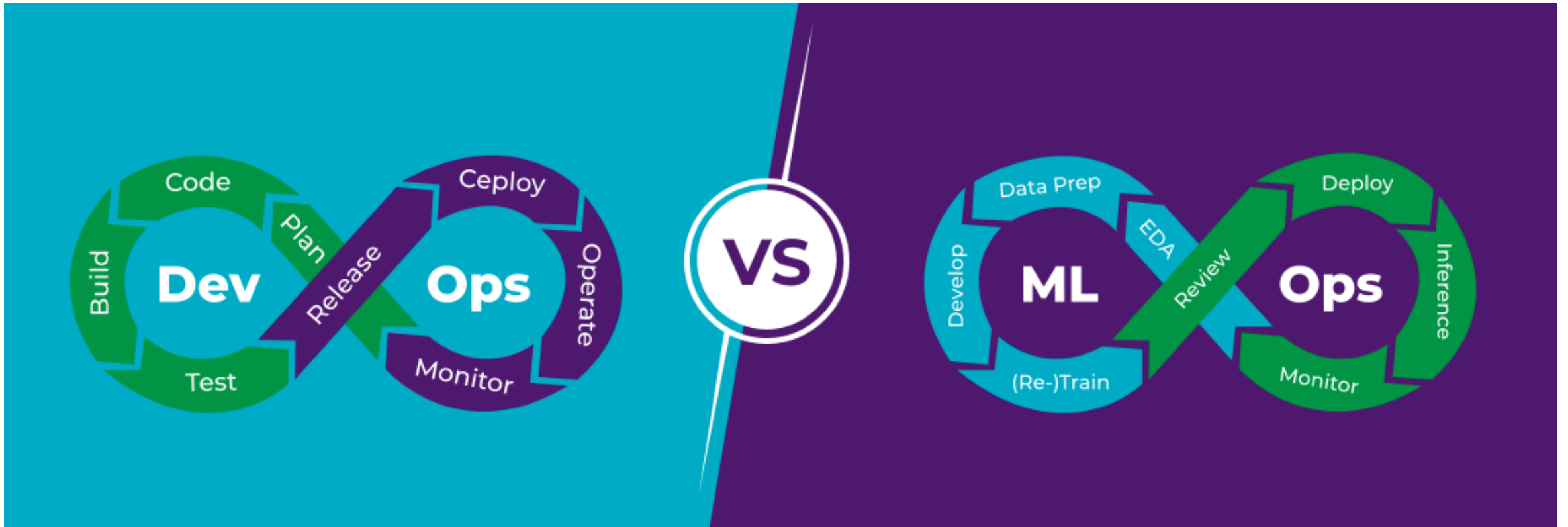
部署方式：

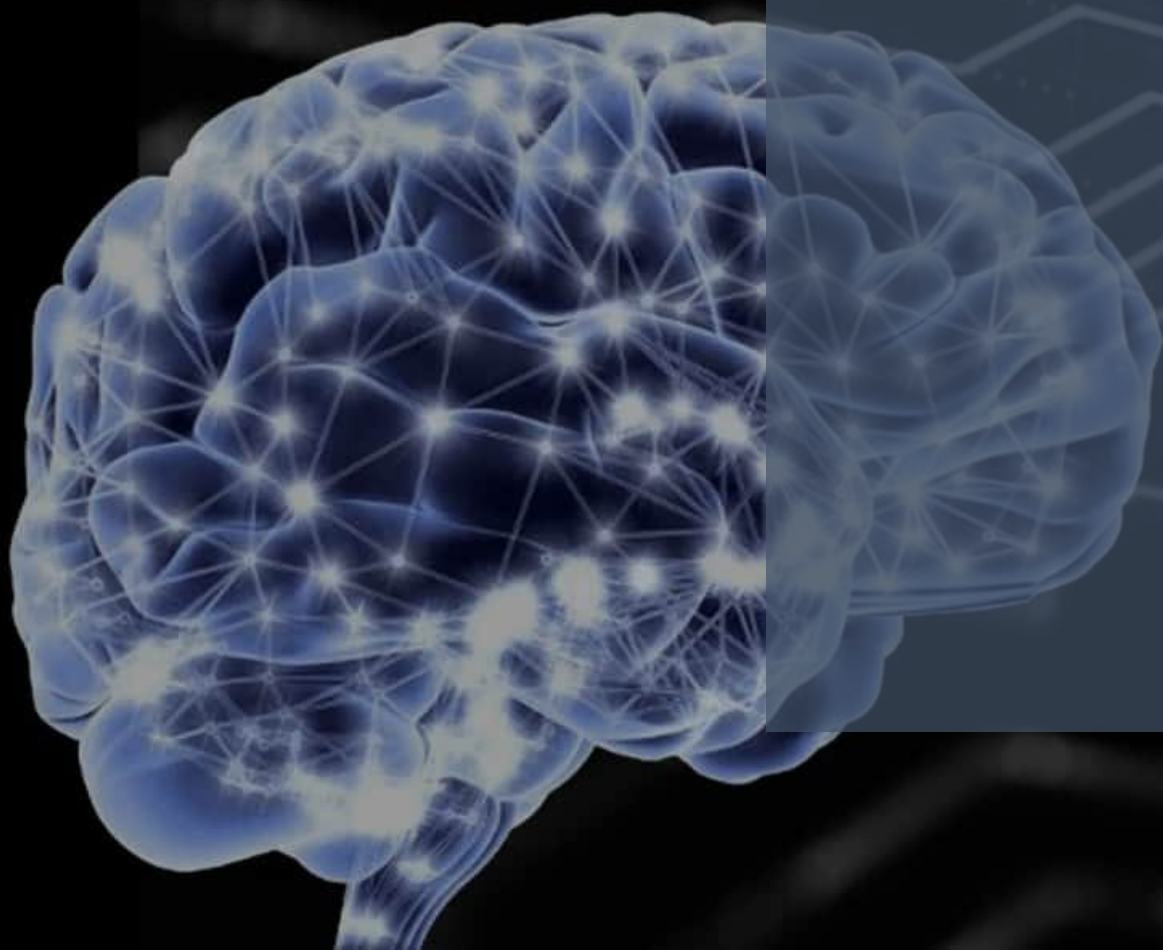
- 训练框架部署引擎：
TensorFlow Lite、TorchServe等；
 - 比完整的训练框架精简；
 - 只支持自己框架训练的模型；
 - 支持的硬件、OS有限。

MLOps



MLOps





谢谢!

文献阅读/专题研讨

- 以小组为单位，进行深度学习领域的**经典文献阅读或专题研讨**
- 小组汇报：
 - PPT介绍（8-10分钟）
 - 答辩、讨论（3-5分钟）
- 每组选择1个主题，并在以下文档登记（每个小组选择的主题不能重复，先到先得）：
<https://docs.qq.com/sheet/DVfVzY2NMqXFSS2FI?tab=BB08J2>
- 评分要点：
 - 结构清晰、内容详略得当、回答问题准确
 - 注意控制时间（时间不足或超时酌情扣分）

文献阅读/专题研讨

日期	周次	小组数量	小组编号
2024.05.27 (周一)	第13周	4个小组	
2024.06.03 (周一)	第14周	4个小组	
2024.06.05 (周三)	第14周	4个小组	
2024.06.17 (周一)	第16周	4个小组	
2024.06.19 (周三)	第16周	5个小组	