

# TensorFlow 遇上 Kubernetes

## 中兴通讯人工智能计算平台的技术实践

刘光聪<sup>1</sup> 韩炳涛<sup>2</sup>

<sup>1</sup>liu.guangcong@zte.com.cn

<sup>2</sup>han.bingtao@zte.com.cn

2017-08-05

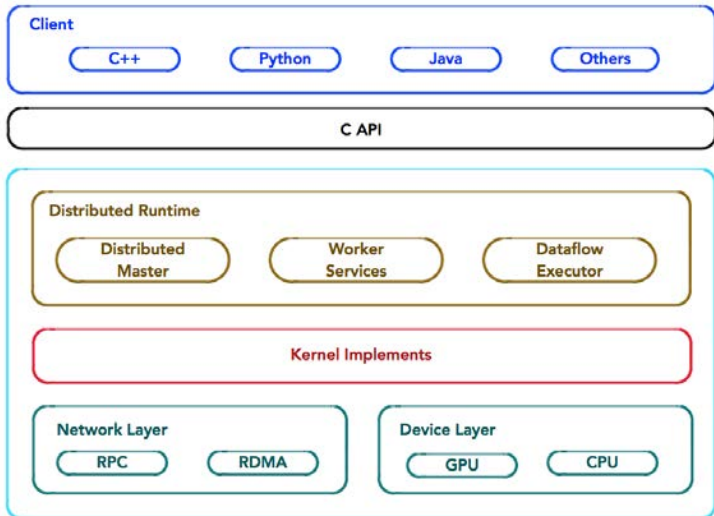
# 内容

- 1 架构概述
- 2 编程模型
- 3 运行模型
- 4 训练模型
- 5 Tensorflow 遇上 Kubernetes
- 6 参考文献

# 架构概述

- 1 系统架构
- 2 图控制

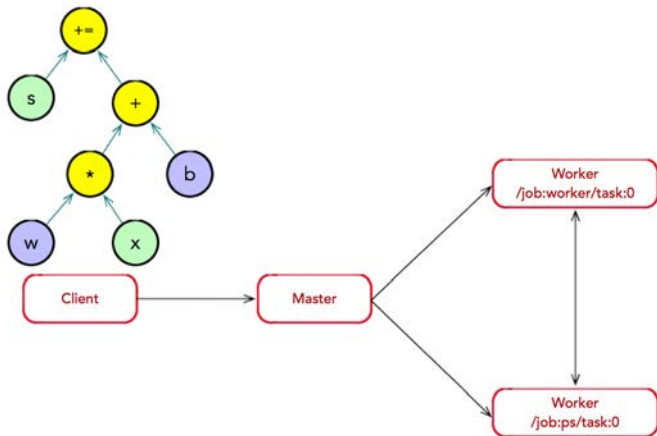
# 系统架构



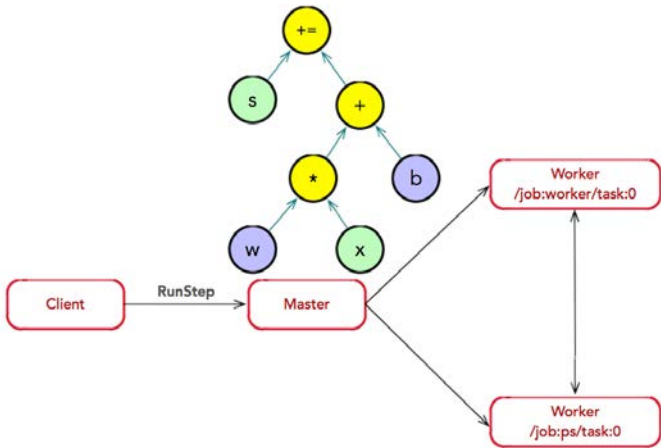
# 设计原则

- **延迟计算**：图的构造与执行分离，并推迟计算图的执行过程
- **原子 OP**：OP 是最小的抽象计算单元，支持构造复杂的网络模型
- **抽象设备**：支持 CPU, GPU, ASIC 多种异构设备类型
- **抽象任务**：基于 Task 的 PS，支持优化算法和网络模型的扩展

## 构造计算图

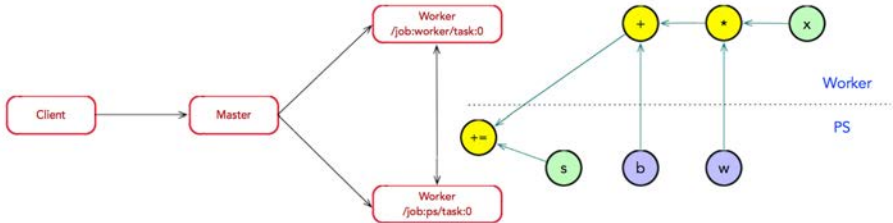


# 执行计算图



图控制

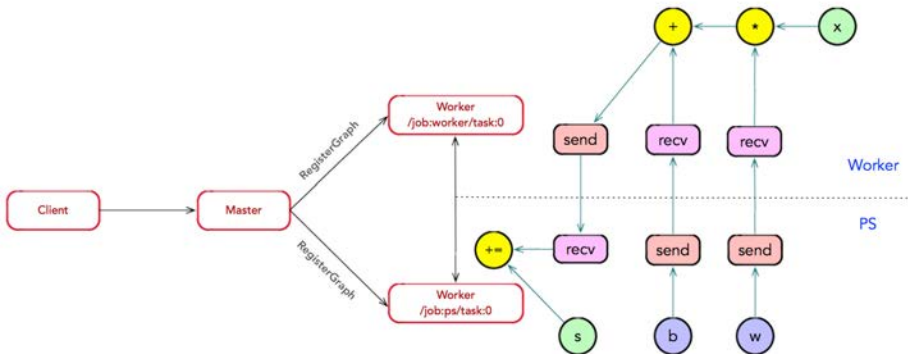
# 图分解：按 Task 分解





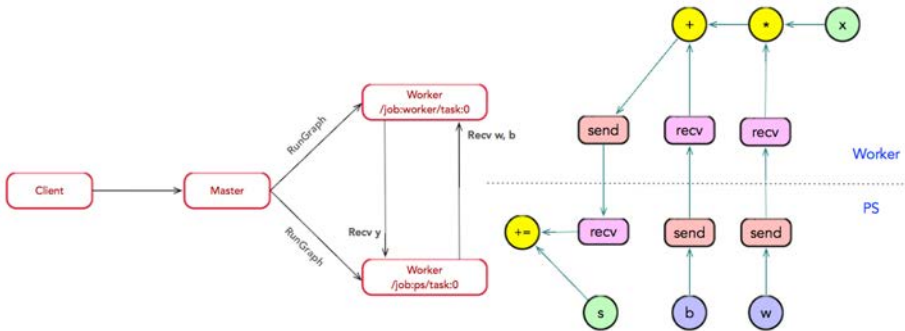
图控制

# 注册图



图控制

# 执行图

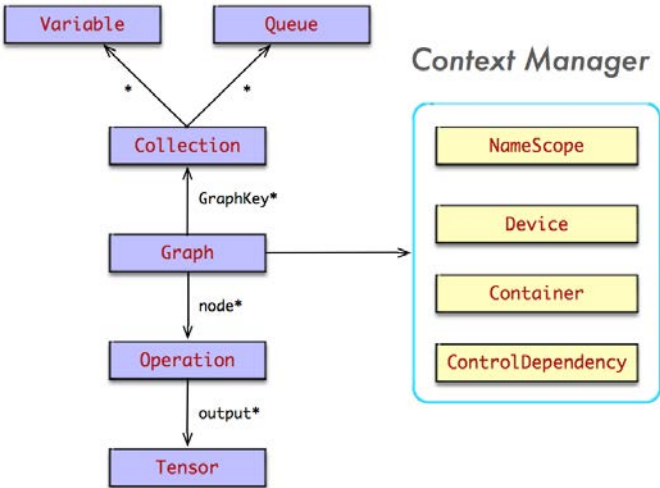


# 编程模型

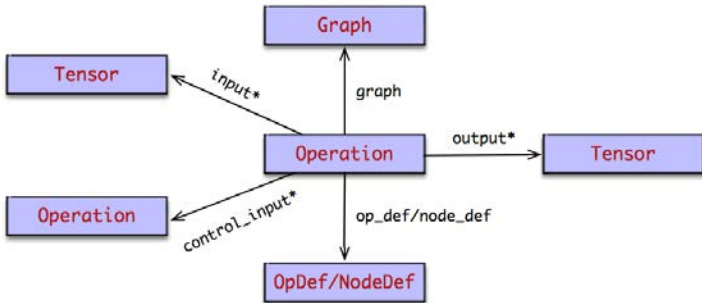
- 1 计算图
- 2 变量
- 3 会话
- 4 图构造

计算图

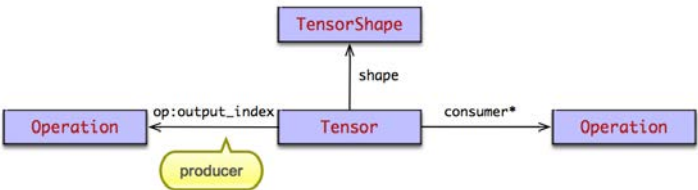
$$Graph = Set\{OP\} + Set\{Tensor\}$$



# OP: 抽象计算

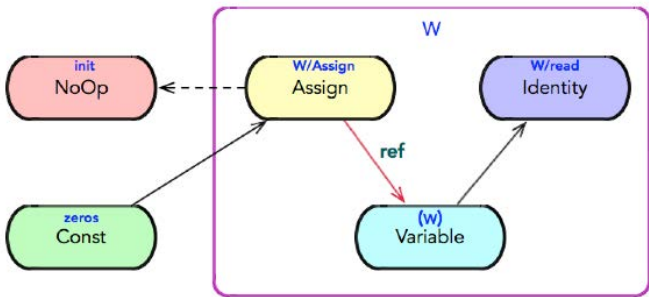


# Tensor: 承载数据



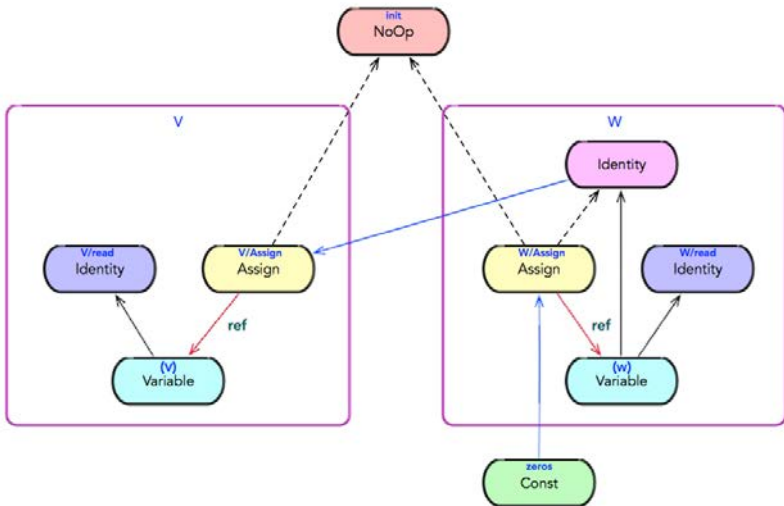
变量

# 初始化模型



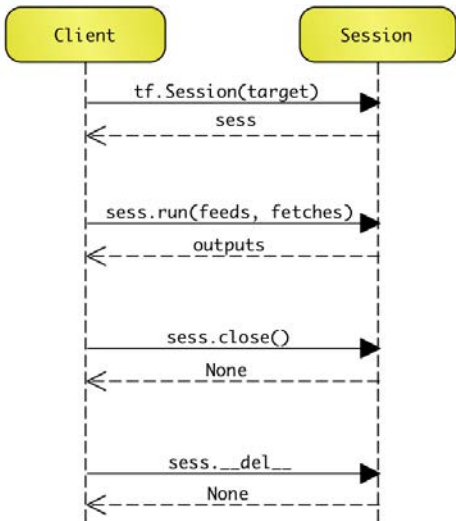
变量

# 初始化依赖



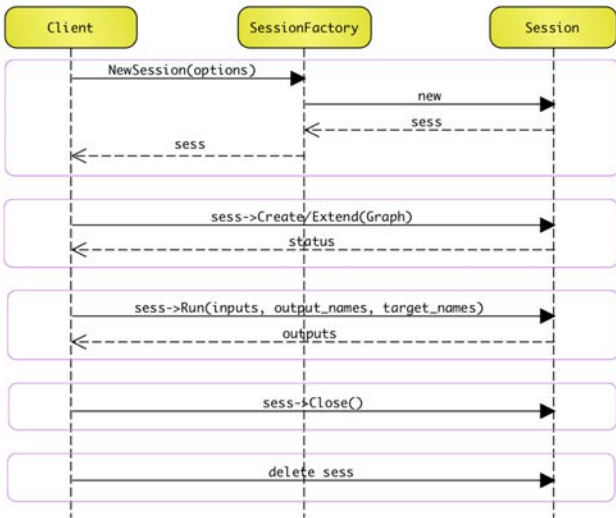


# 生命周期: Python

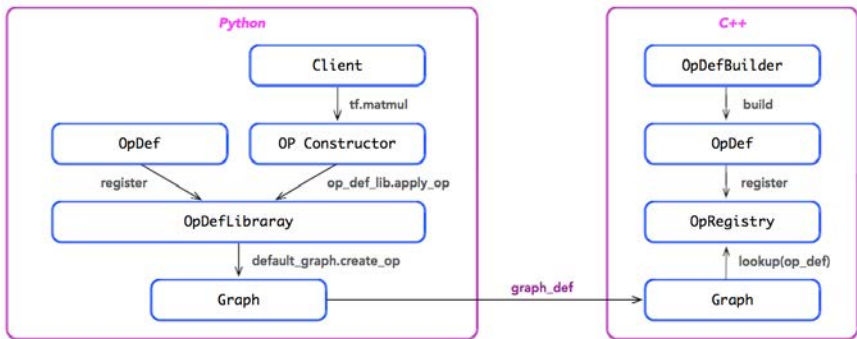


会话

# 生命周期: C++



# 图构造与传递



# 实例: OP 构造器

## OP Constructor

```
def zeros_like(tensor, name=None):
    gen_array_ops._zeros_like(tensor, name=name)

tensor = tf.constant([1, 2], name="n1")
zeros = tf.zeros_like(tensor, name="n2")
```

## Code Generator

```
def _zeros_like(dtype, shape=None, name=None):
    return _op_def_lib.apply_op("ZerosLike", x=x, name=name)
```

# 实例：构造 OP

## OpDef Repository

```
class OpDefLibrary(object):
    def apply_op(self, op_name, name=None, **keywords):
        inputs, input_types, output_types, attr_protos, op_def =
        with graph.as_default(), ops.name_scope(name) as scope:
            return graph.create_op(op_name, inputs, output_types, name=scope,
                                   input_types=input_types, attrs=attr_protos, op_def=op_def)
```

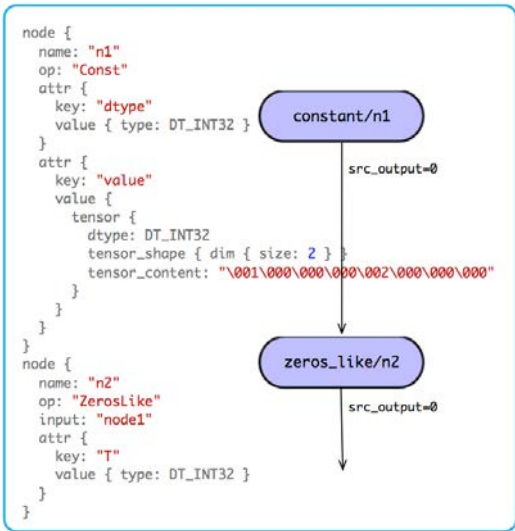


## Graph

```
class Graph(object):
    def create_op(self, op_type, inputs, dtypes, input_types=None,
                 name=None, attrs=None, op_def=None):
        node_def, control_inputs = ...
        return Operation(node_def, self, inputs=inputs, output_types=output_types,
                        control_inputs=control_inputs, input_types=input_types,
                        op_def=op_def)
```



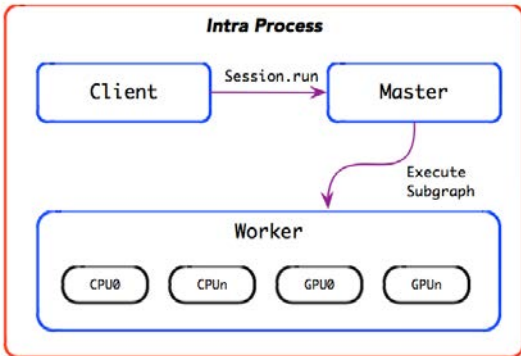
# 实例：图构造



# 运行模型

- 1 运行模式
- 2 创建会话
- 3 迭代执行

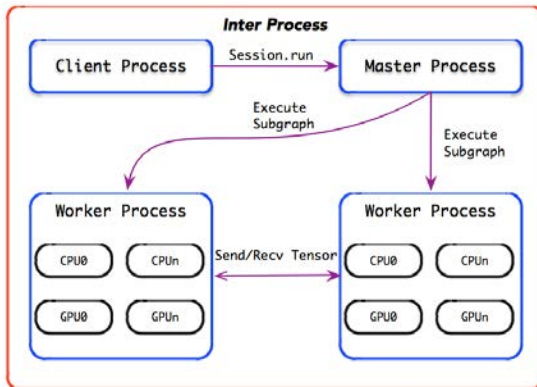
# 本地模式



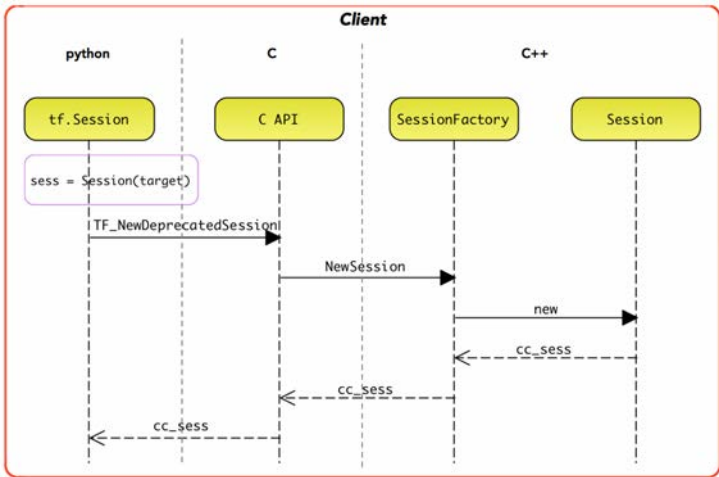


运行模式

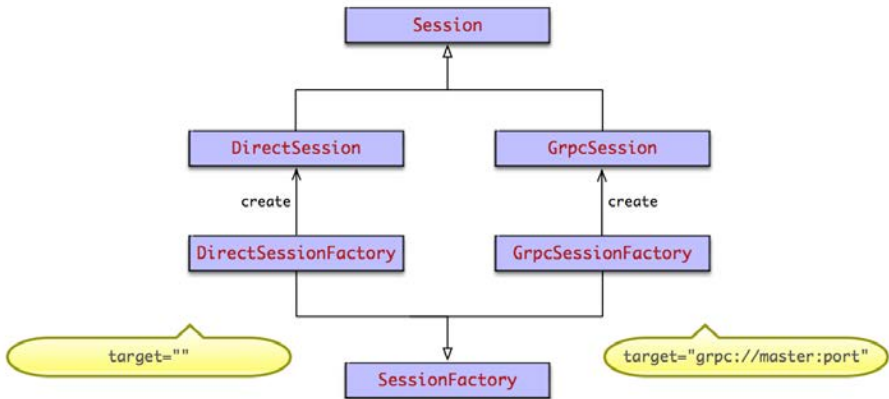
## 分布式模式



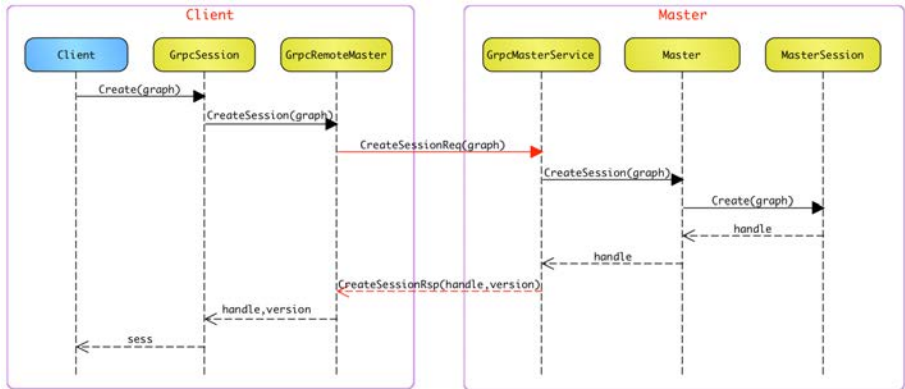
# 创建 ClientSession



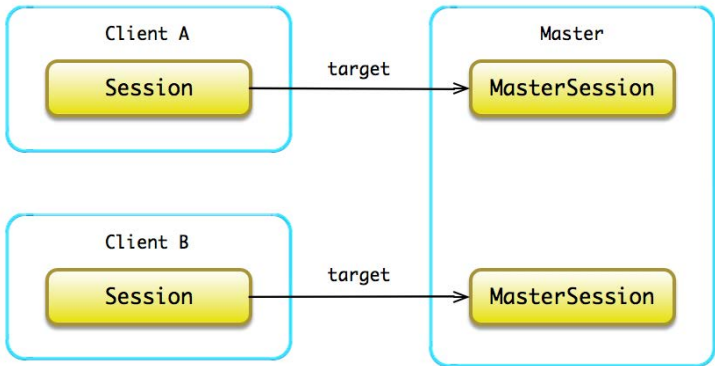
# 多态创建



# 创建 MasterSession

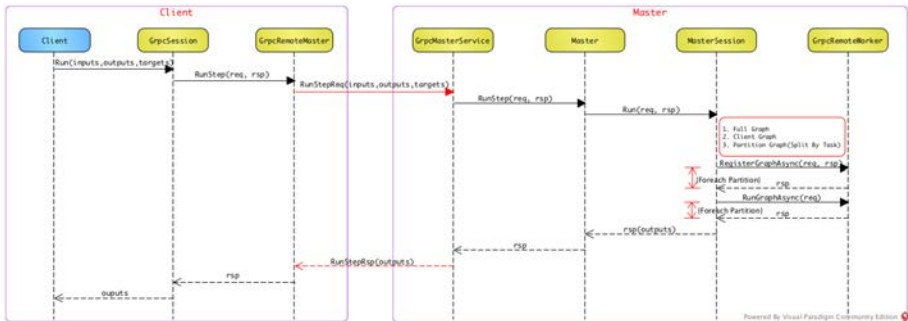


# MasterSession 模型



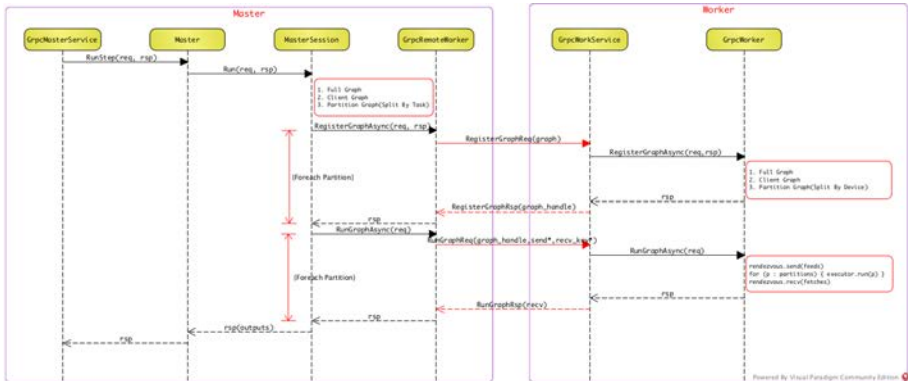
迭代执行

# 一级图分裂



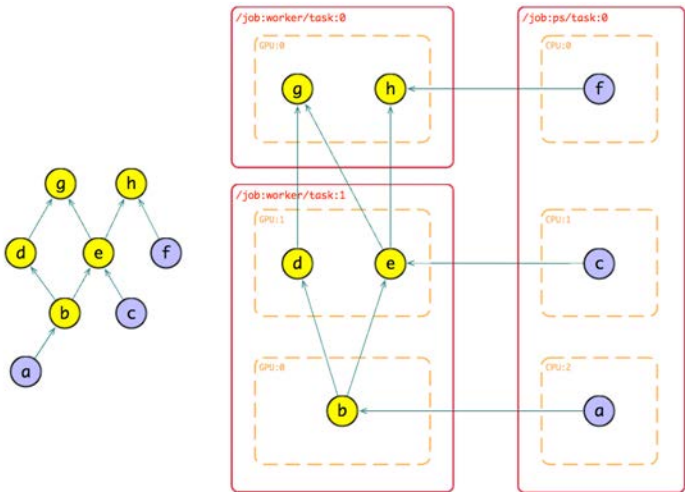
迭代执行

# 二级图分裂



迭代执行

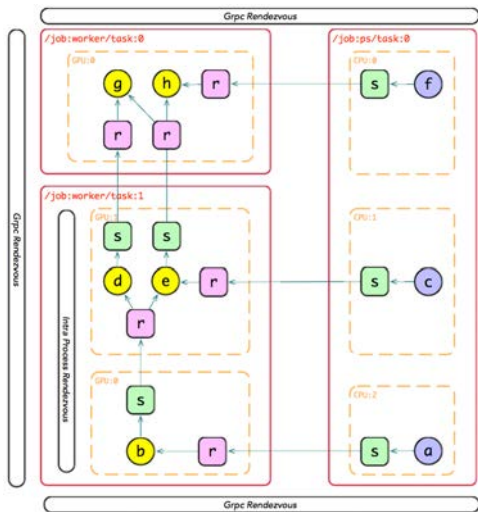
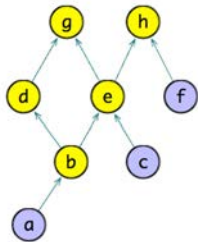
# 实例：图分裂





迭代执行

## 实例：交换数据



# 训练模型

- 1 计算梯度
- 2 应用梯度
- 3 工作流

# 优化器: 计算梯度

```

class Optimizer(object):
    def minimize(self, loss, var_list=None, global_step=None):
        grads_and_vars = self.compute_gradients(
            loss, var_list=var_list)
        return self.apply_gradients(
            grads_and_vars,
            global_step=global_step)

    def compute_gradients(loss, var_list):
        grads = gradients(loss, var_list, grad)
        return list(zip(grads, var_list))

    def gradients(loss, var_list, grads=1):
        ops_and_grads = {}
        for op in reversed_graph(loss).topological_sort():
            grad = op.grad_fn(grad)
            ops_and_grads[op] = grad
        return [ops_and_grads.get(var) for var in var_list]

```

# 梯度函数

```
@ops.RegisterGradient("op_name")
def grad_func(op, grad):
    """construct gradient subgraph for an op type.
    Returns:
        A list of gradients, one per each input of op.
    """
    return cons_grad_subgraph(op, grad)
```

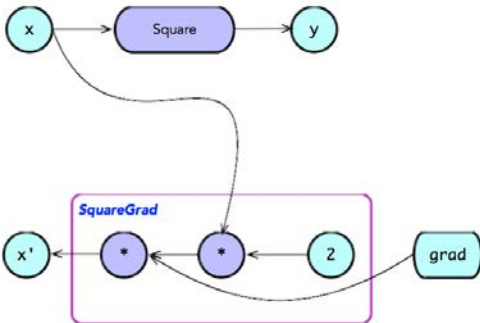
$$(y_1, y_2, \dots, y_m) = f(x_1, x_2, \dots, x_n)$$

$$(\partial L / \partial x_1, \partial L / \partial x_2, \dots, \partial L / \partial x_n) = g(x_1, x_2, \dots, x_n; \partial L / \partial y_1, \partial L / \partial y_2, \dots, \partial L / \partial y_m)$$

# 例子：SquareGrad 函数

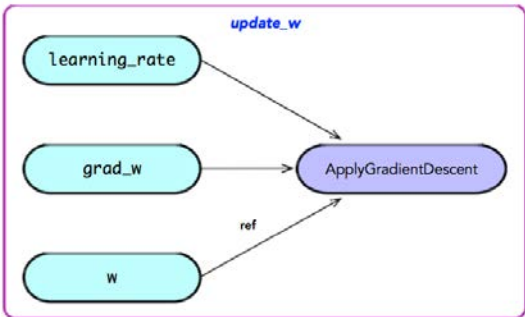
```

@ops.RegisterGradient("Square")
def SquareGrad(op, grad):
    x = op.inputs[0]
    with ops.control_dependencies([grad.op]):
        return grad * (2.0 * x)
    
```

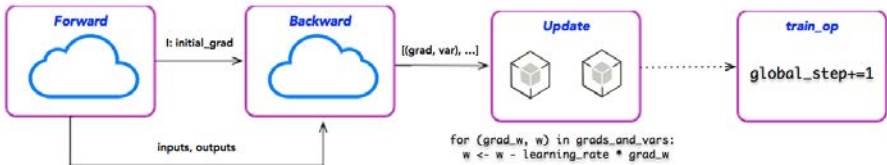


# 应用梯度

```
def apply_gradients(grads_and_vars, learning_rate):
    for (grad, var) in grads_and_vars:
        apply_gradient_descent(learning_rate, grad, var)
```

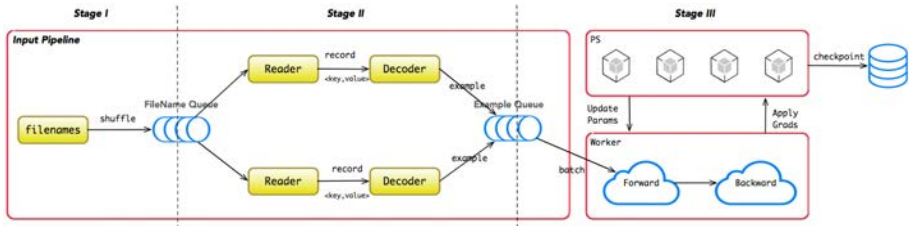


# RunStep 过程



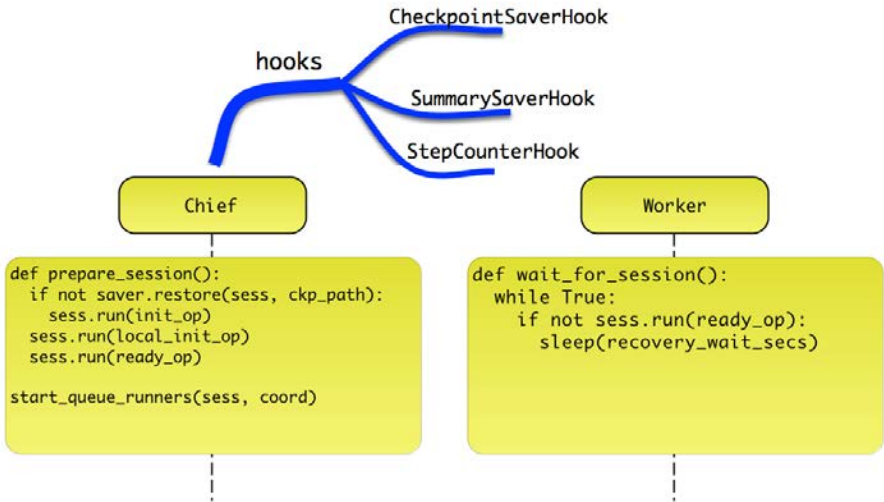
工作流

# 大规模模型训练





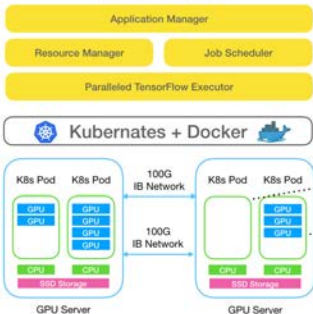
# 初始化模型



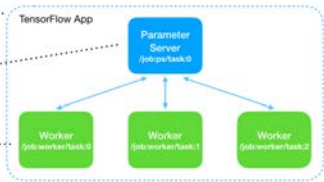
# Tensorflow 遇上 Kubernetes

# 容器化高性能计算集群

## 云计算平台的架构



## 分布式TF程序在K8s上的调度



- 使用Nvidia DGX-1搭建高性能计算集群
- 基于Docker容器使用GPU资源，在一个容器中映射多个GPU设备
- 使用Kubernetes管理集群资源，为每个训练/推理作业创建多个POD

- Application Manager: 管理每个TensorFlow App的Pod List
- Resource Manager: 决定为一个App分配哪些Pod
- Job Scheduler: 决定在Pod上部署哪些TF进程(Worker or PS)
- Executor: 运行并管理容器内的TF进程

# 扩展 K8s 支持 GPU 和 RDMA 网络

原生K8s不支持GPU和RDMA网络资源的管理，我们对K8s进行了如下扩展：

- 能够从K8s查询节点的GPU设备信息，如GPU个数、显存大小等
- 能够在Pod List中指定映射哪个GPU
- 能够查询节点RDMA网络的信息
- 支持-cpuset-cpus选项，在某些场景能够提升100%性能。
- 资源分配策略是最容易变化的部分。为了减少对K8s的影响，我们把K8s看做是执行器而不是资源调度器。分离了单独的RM子系统来实现多种调度策略（为了简化RM，假设K8s集群只运行TF业务）。

```
apiVersion: v1kind: Pod
metadata:
  name: s1
spec:
  containers:
  - name: s1
    image: library/nginx:v1.0
    imagePullPolicy: IfNotPresent
  resources:
    limits:
      alpha.kubernetes.io/gpuset: 0,1,3
  nodeSelector:
    kubernetes.io/hostname: 10.114.51.204
```

在Pod中指定GPU

```
"labels": {
  "beta.kubernetes.io/arch": "amd64",
  "beta.kubernetes.io/os": "linux",
  "kubernetes.io/hostname": "10.114.51.204",
  "alpha.kubernetes.io/infiniBand": "true"
}
```

RDMA能力发现，通过给相应的节点打标签方式实现

Tensorflow On Kubernetes

# 在云上执行分布式 TF 计算

## 在云环境运行TF并不轻松

- 云环境下每次分配的资源是动态的

```

# Keras on AWS EC2
$ python trainer.py
--k8s-namespace aws-1022-p01-aws-ec2-us-1022-1
--worker_aws-ec2-us-1022-p01-aws-ec2-us-1022-1
--job_name keras --task train

# Keras on Azure VM
$ python trainer.py
--k8s-namespace aws-1022-p01-aws-ec2-us-1022-1
--worker_aws-ec2-us-1022-p01-aws-ec2-us-1022-1
--job_name keras --task train

# Keras on GCP VM
$ python trainer.py
--k8s-namespace aws-1022-p01-aws-ec2-us-1022-1
--worker_aws-ec2-us-1022-p01-aws-ec2-us-1022-1
--job_name keras --task train
    
```

- 硬编码分布式方案不适合动态的环境（希望弹性计算）

```

with tf.device('/job:ps/task:0'):
    weights_1 = tf.Variable(...)
    biases_1 = tf.Variable(...)

with tf.device('/job:ps/task:1'):
    weights_2 = tf.Variable(...)
    biases_2 = tf.Variable(...)
    
```

不灵活，写起来也很烦！

- 希望最大化的优化性能（怎样分布式是最合理的）

## 解决之道：自动化实现TF的分布式执行



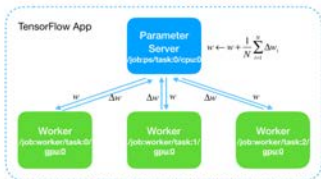
TensorFlow App 自动化分布式执行流程



机器学习算法预测+模拟退火算法优化，相对于简单贪心策略提升20%-25%的性能

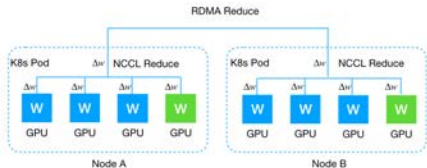
# 计算模型优化

## 经典的分布式计算模型效率低下



- 跨Pod之间的通信需要走容器网络，带来额外开销
- GPU-CPU之间的数据传输非常慢

## 全GPU计算模型



- 如果TF app部署在一个节点，只创建一个容器，映射多个GPU。利用NCCL all-reduce在GPU之间完成参数更新，不经过CPU，也不使用Parameter Server。
- 如果TF app部署跨多个节点，则在每个节点上创建一个容器。每个节点内，先将参数更新one-reduce到一张GPU卡上（主卡），再在每个节点的主卡之间用RDMA网络执行all-reduce。最后主卡将更新后的参数利用NCCL broadcast下发给本节点的所有GPU。



# 文献

- TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems, Google Inc.
- TensorFlow: A System for Large-Scale Machine Learning, Google Inc.



致谢

# Q&A



致谢

# 致谢

# Thanks