

# 基于 k8s 原生扩展的机器学习 平台引擎 ML Engine

褚向阳

小米 高级软件工程师





# SPEAKER INTRODUCE

---

## 褚向阳 高级软件工程师

- 2013 年毕业后加入红帽软件，吸收开源文化，接触 OpenStack 和 IaaS 平台相关技术。
- 2015 年底开始加入容器云创业公司数人云，参与打造容器化的 PaaS 平台。
- 2018 年从京东广告部加入小米人工智能部，负责小米机器学习平台的建设，重点支持各个框架的分布式训练，订制优化 K8s 调度，努力提高平台用户体验的同时保证集群利用率。
- 持续关注 Kubeflow 社区及性能优化相关开源项目发展。





# TABLE OF CONTENTS 大纲

---

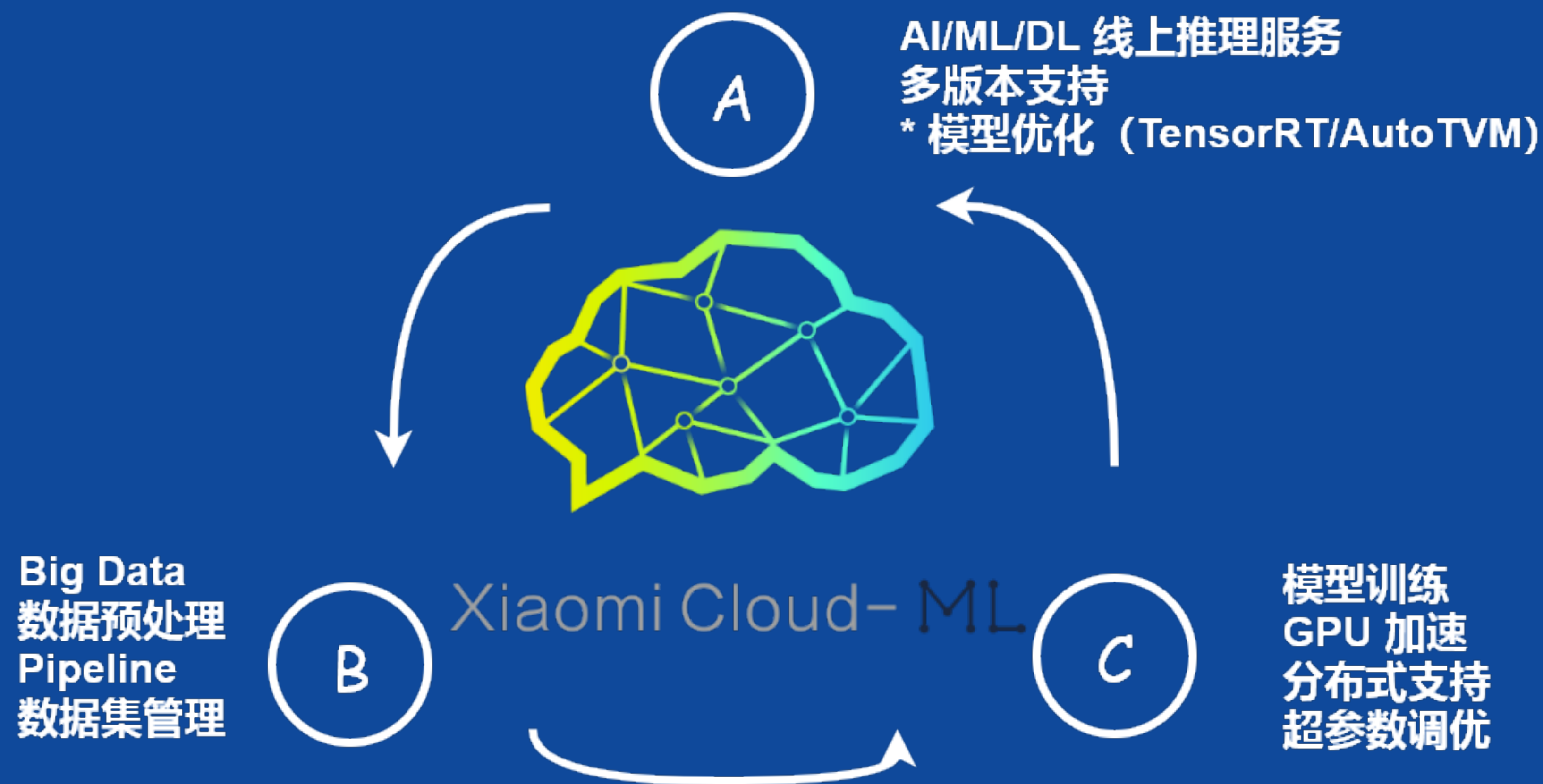
- 小米 CloudML 机器学习平台简介
- ML Engine 架构设计演进
- ML Engine 对多框架的分布式训练支持详解
- 未来发展方向和具体工作

# TABLE OF CONTENTS 大纲

---

- 小米 CloudML 机器学习平台简介
  - ML Engine 架构设计演进
  - ML Engine 对多框架的分布式训练支持详解
  - 未来发展方向和具体工作

# CloudML & ABC



# 平台全景





# 平台全景

核心  
逻辑



ML  
Engine

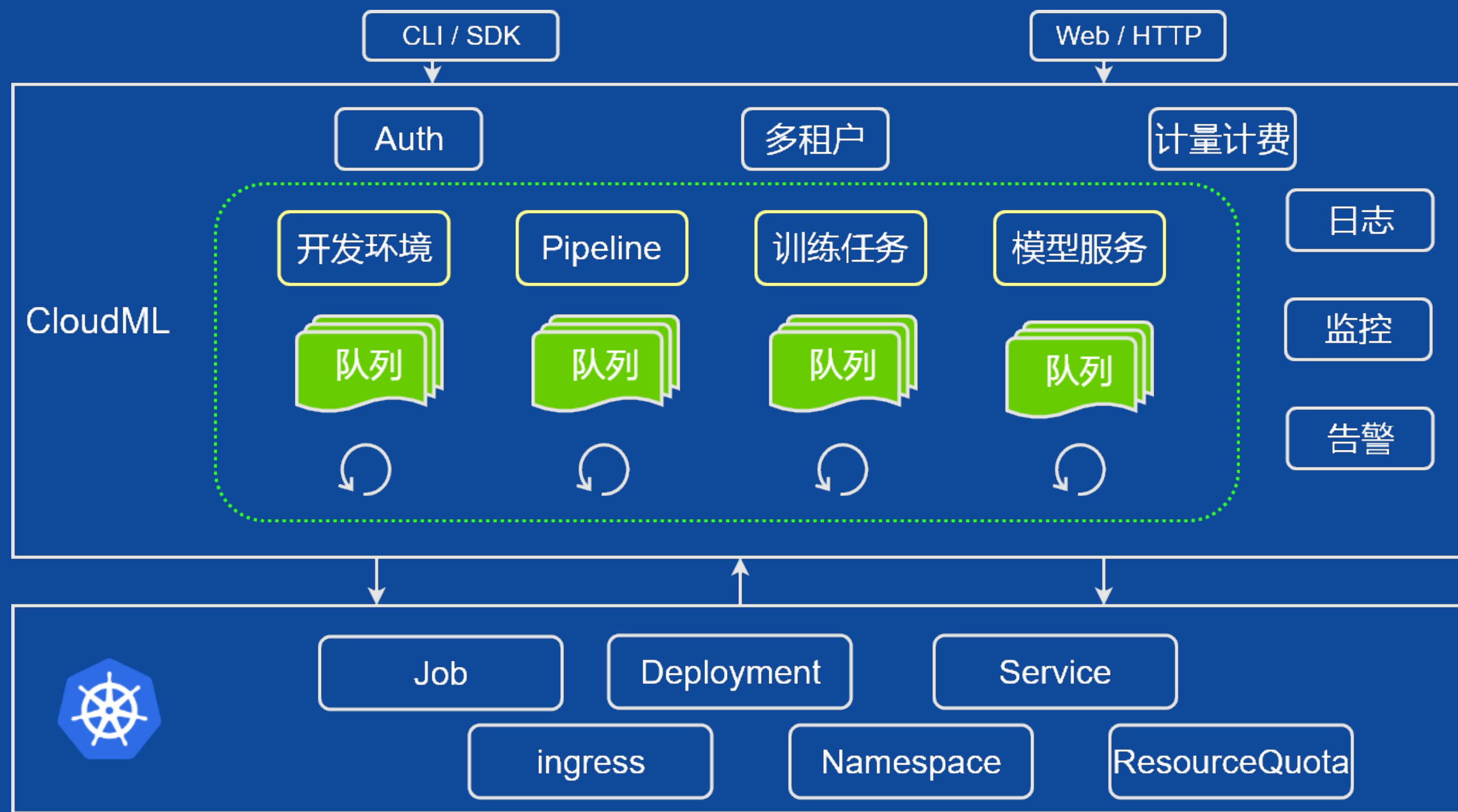
# TABLE OF CONTENTS 大纲

---

- 小米 CloudML 机器学习平台简介
- **ML Engine 架构设计演进**
- ML Engine 对多框架的分布式训练支持详解
- 未来发展方向和具体工作



# CloudML 架构



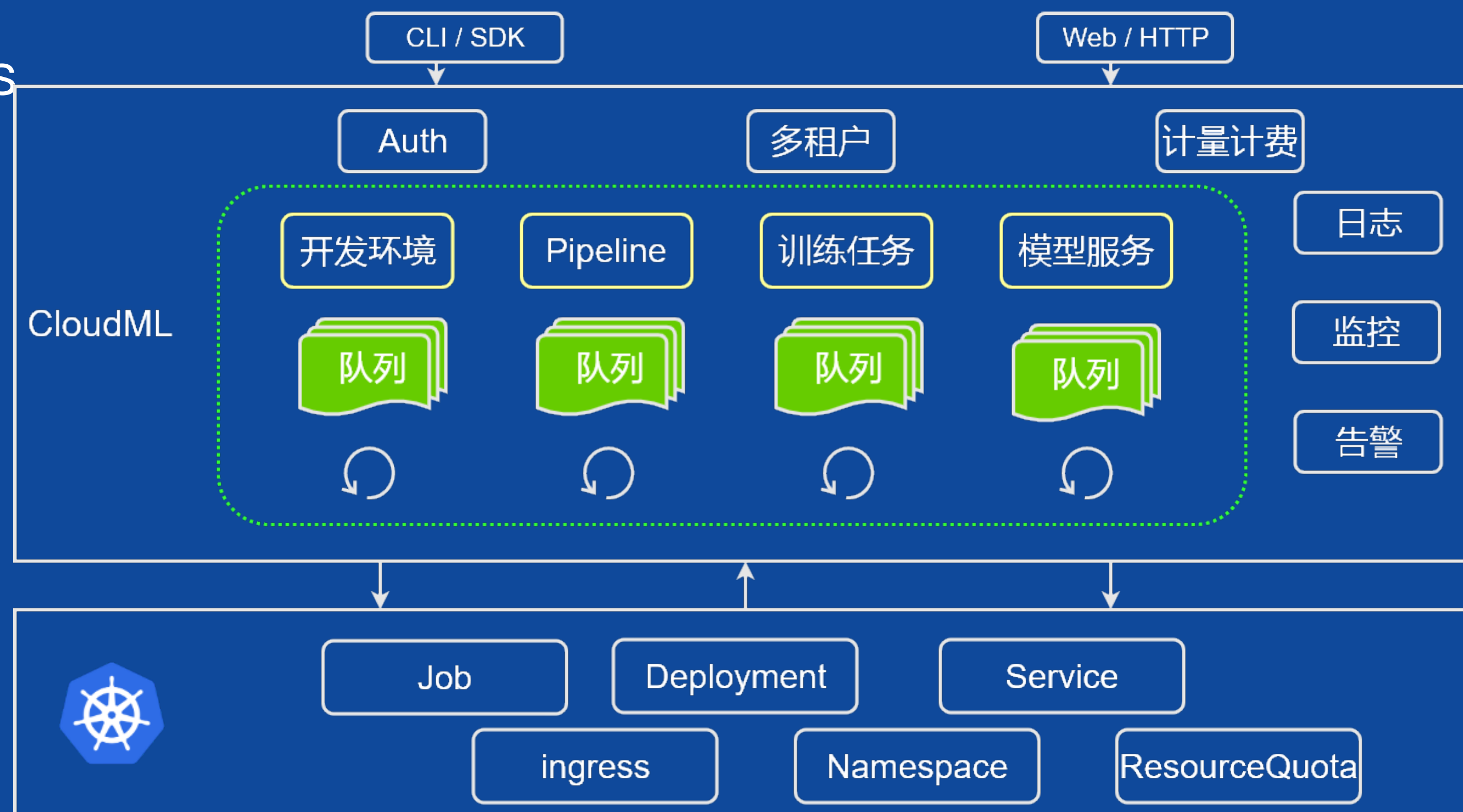
# 架构和挑战

## 机器学习平台架构

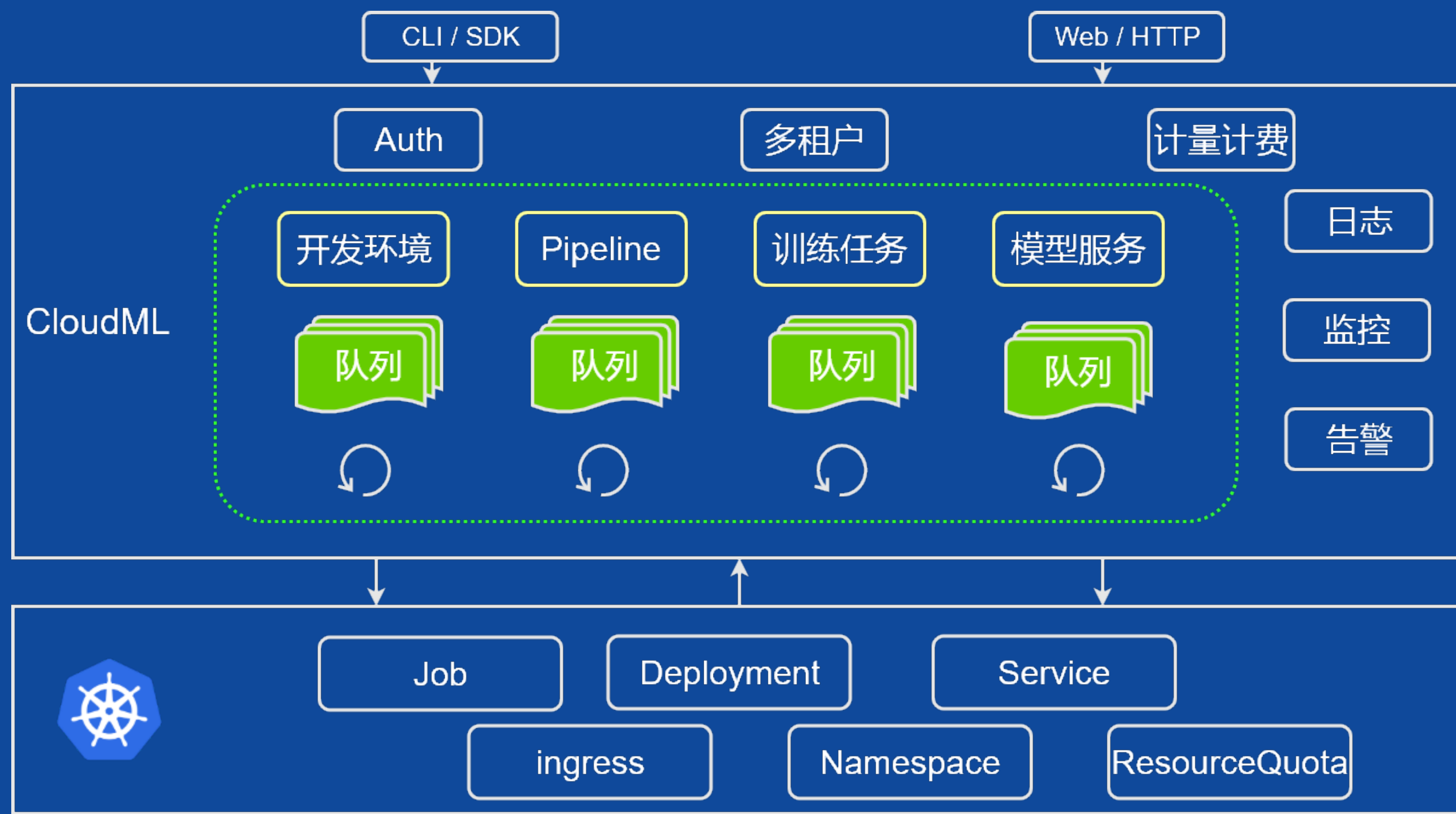
- 分布式训练任务 -> 一组 Job/Pod + SVC
- 模型服务 -> Deployment + SVC + Ingress

## 挑战:

- 队列、状态同步
- 生命周期管理
- 调度策略



# CloudML 架构



ML Engine  
?



The easiest way to install  
Kubernetes is to get someone  
else to do it for you.

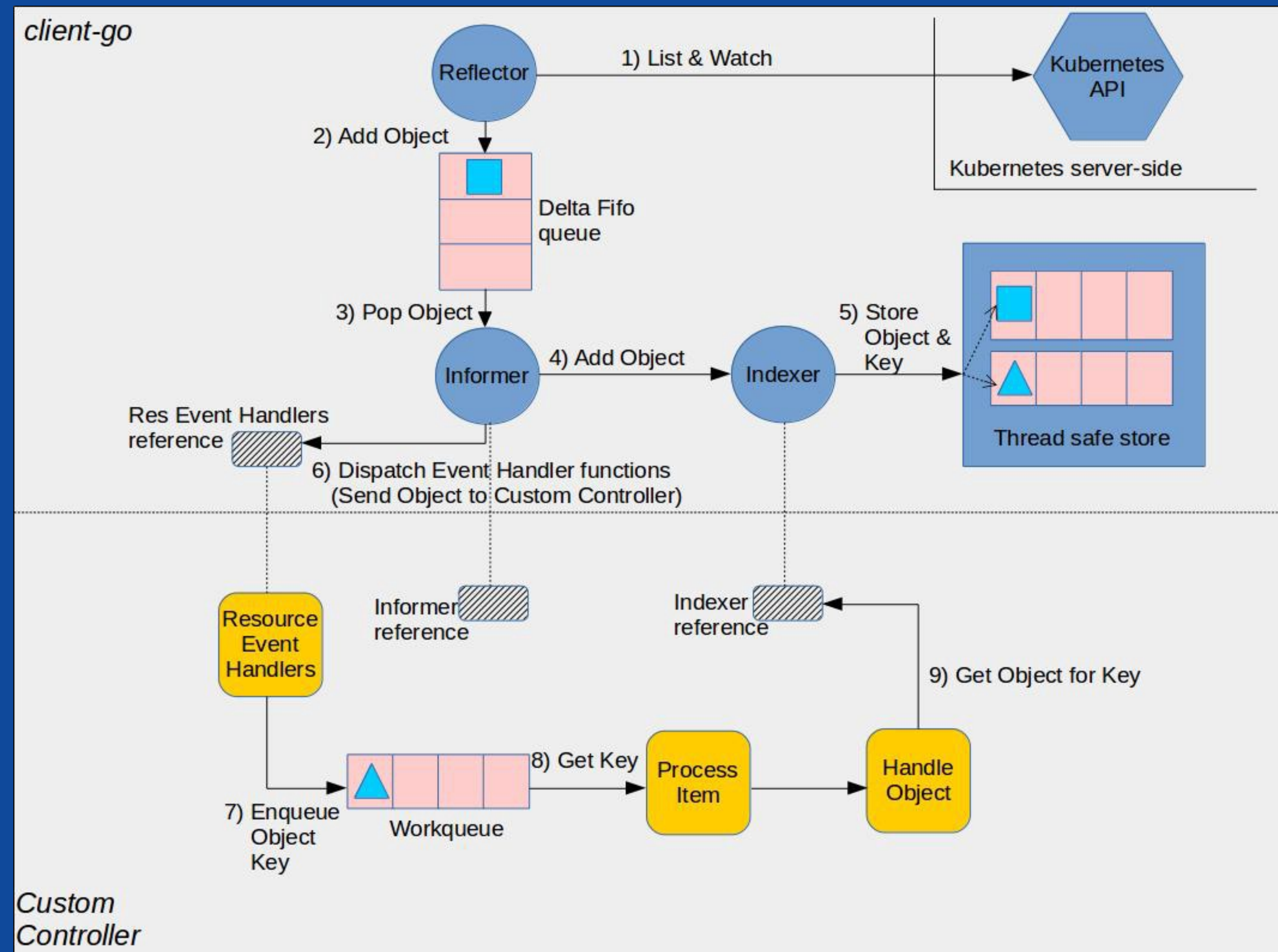
#LISA19, @jpetazzo

The easiest way to Build a  
AI /ML platform is to let k8s  
to do it for you.

#ArchSummit19, me

# K8S 扩展资源定义

## CRD & Controller

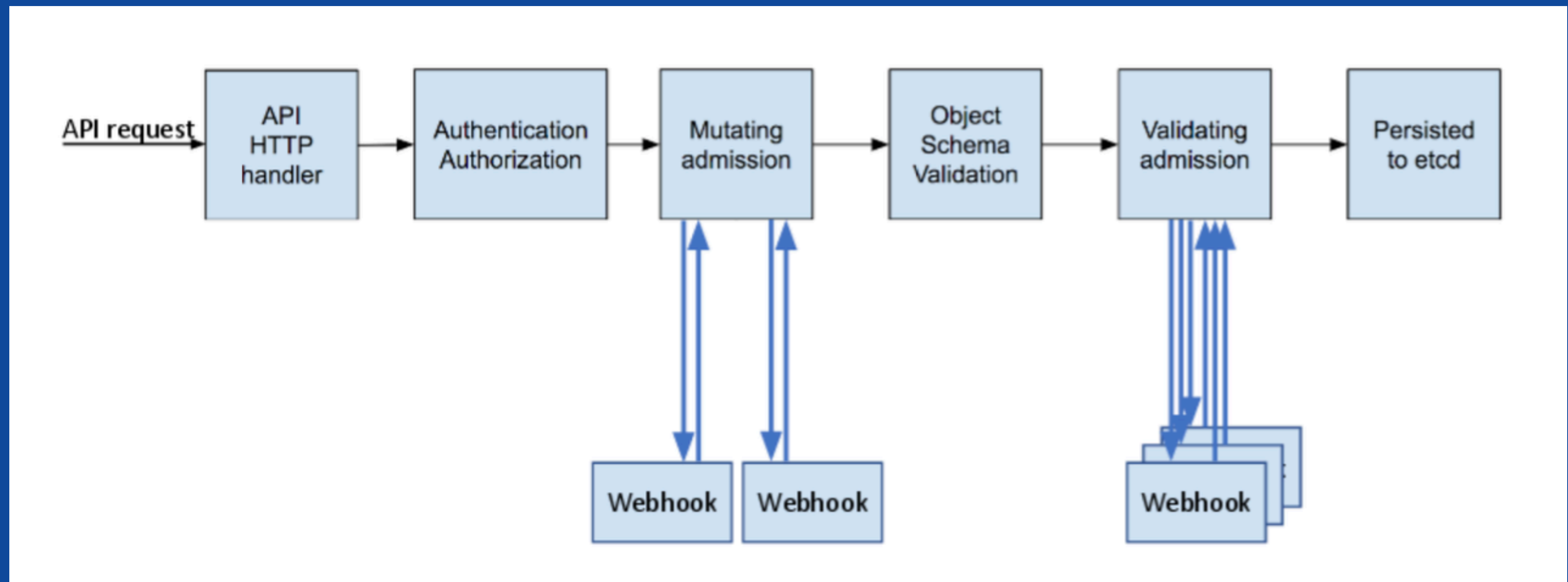


<https://medium.com/@cloudark/kubernetes-custom-controllers-b6c7d0668fdf>



# Mutating & Validating

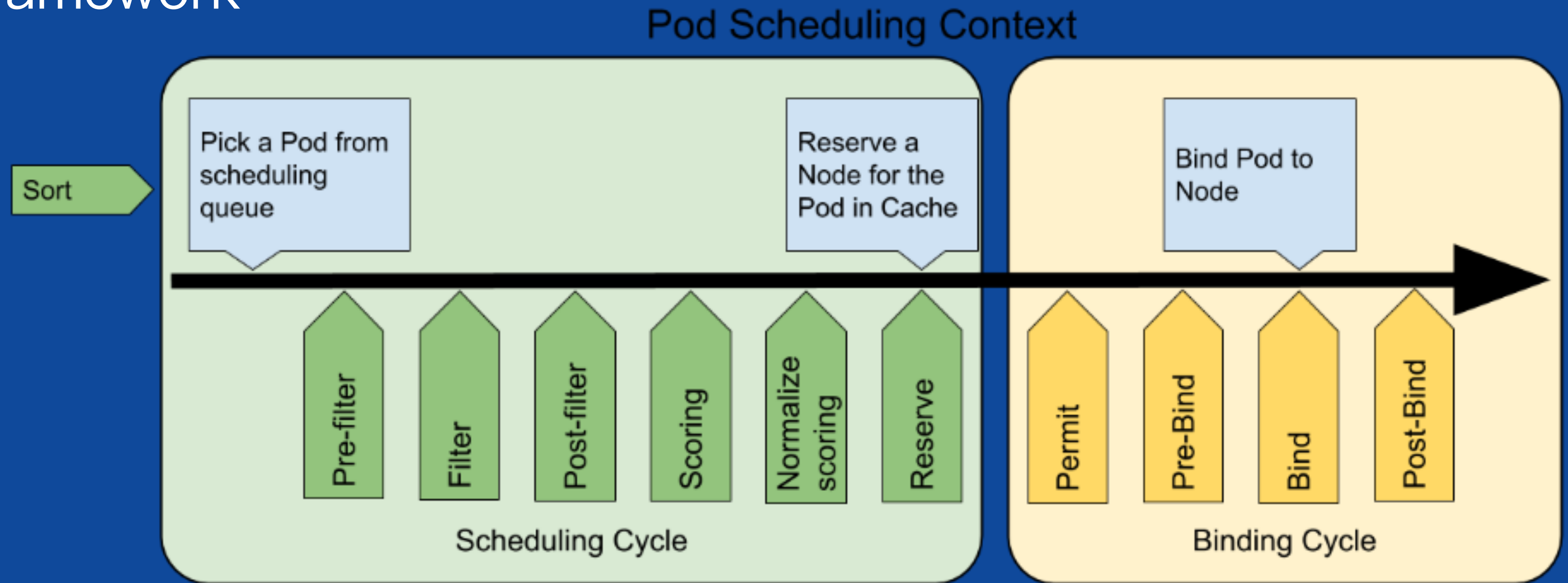
## Admission webhooks



<https://kubernetes.io/blog/2019/03/21/a-guide-to-kubernetes-admission-controllers/>

# 调度优化

## Scheduling framework



<https://kubernetes.io/docs/concepts/configuration/scheduling-framework/>

# 扩展工具包

## Kube-builder

- <https://github.com/kubernetes-sigs/kubebuilder>

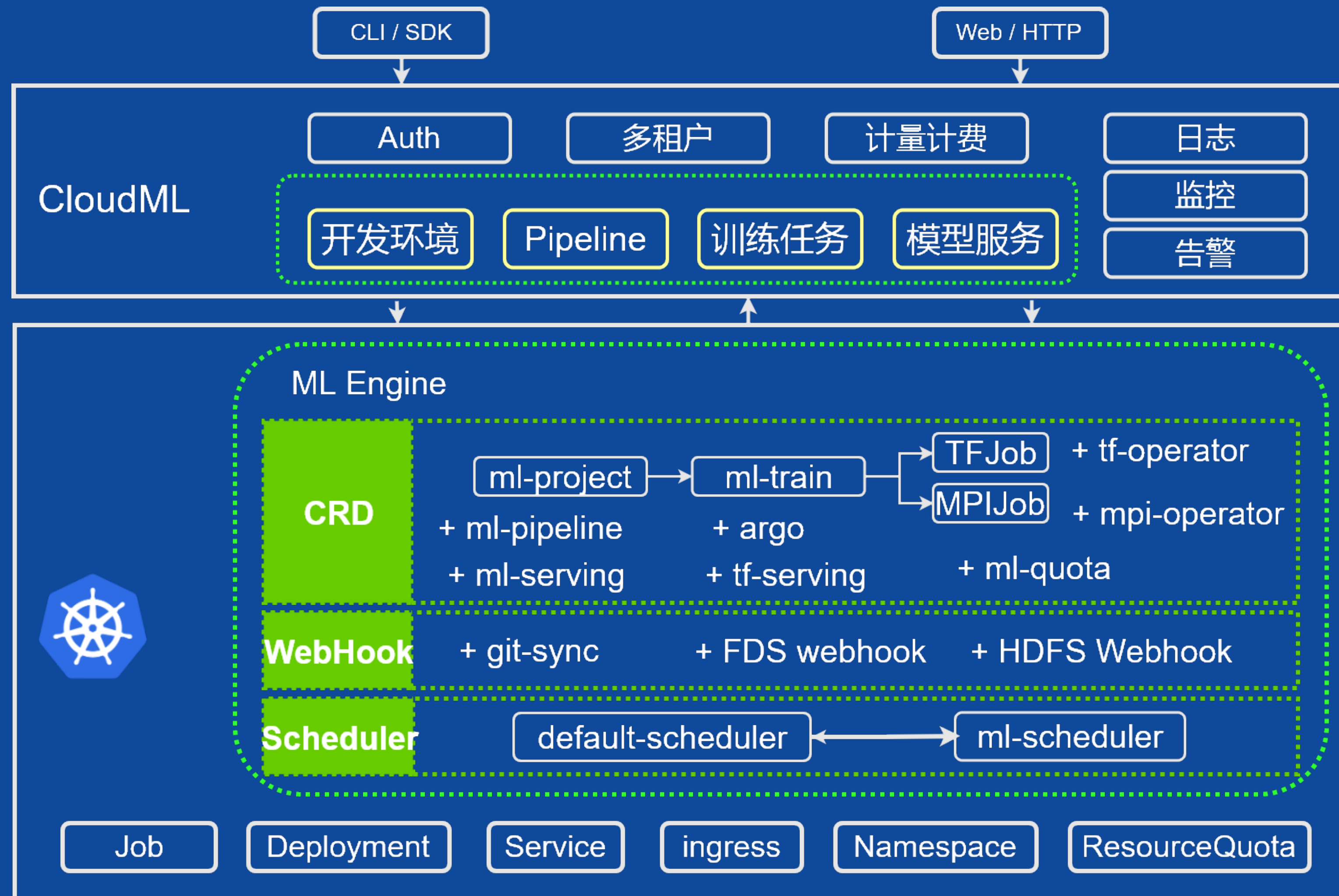
## Operator Framework/Operator SDK

- <https://github.com/operator-framework/operator-sdk>

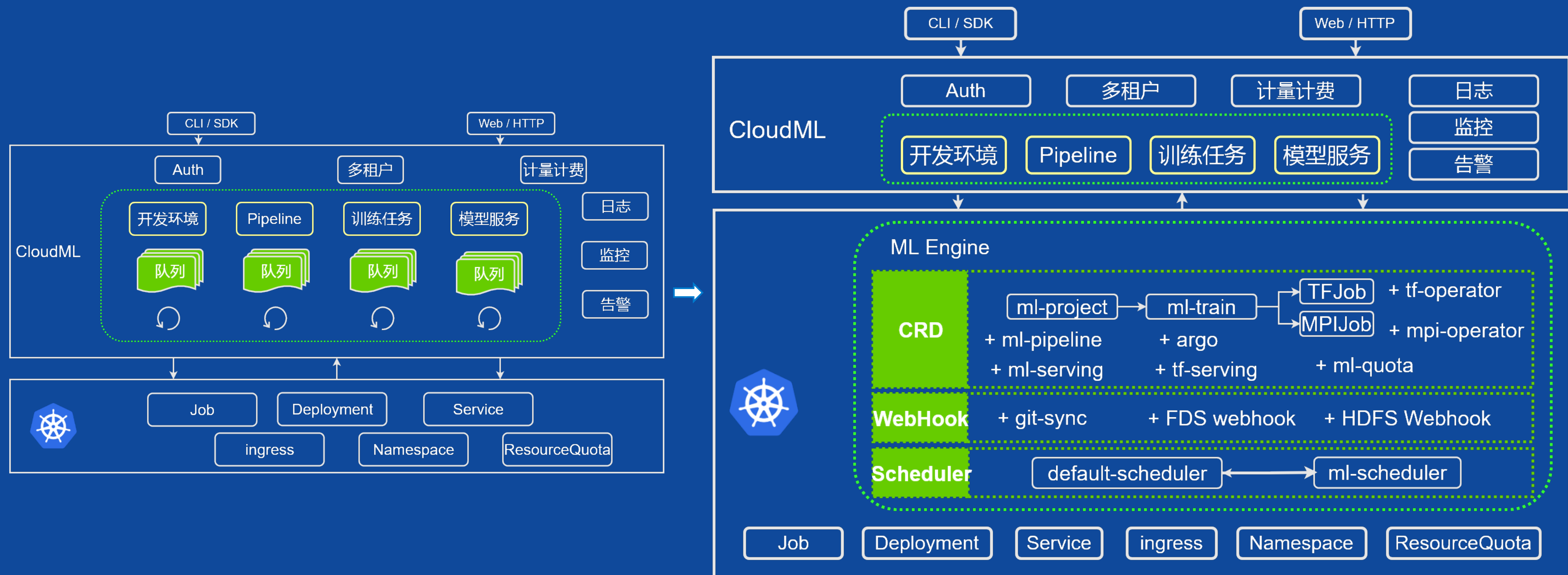
```
$ # Initialize Go modules
$ go mod init demo.kubebuilder.io
go: creating new go.mod: module demo.kubebuilder.io
$ # Let's initialize the project
$ kubebuilder init --domain tutorial.kubebuilder.io
go get sigs.k8s.io/controller-runtime@v0.2.2
go mod tidy
Running make...
make
$HOME/go/bin/controller-gen object:headerFile=./hack/boilerplate.go.txt paths="./..."
go fmt ./...
go vet ./...
go build -o bin/manager main.go
```



# ML Engine 最终架构



# 架构对比



# ML Engine 架构设计演进

## ML Engine 新版架构总结

- 机器学习平台业务的合理抽象 CRD
- 针对机器学习任务特性的定制调度器
- 公共的 validating/mutating 逻辑
- 提供统一的对外接口



# TABLE OF CONTENTS 大纲

---

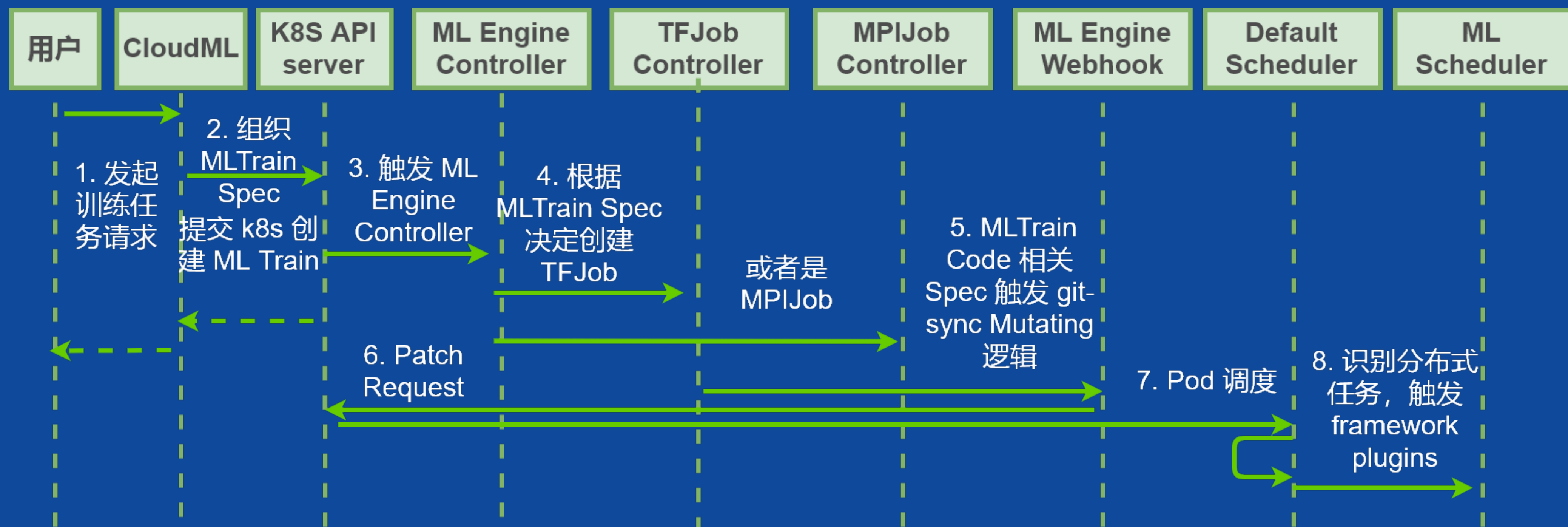
- 小米 CloudML 机器学习平台简介
- ML Engine 架构设计演进
- **ML Engine 对多框架的分布式训练支持详解**
- 未来发展方向和具体工作

# ML Engine 训练任务需求

核心用例：

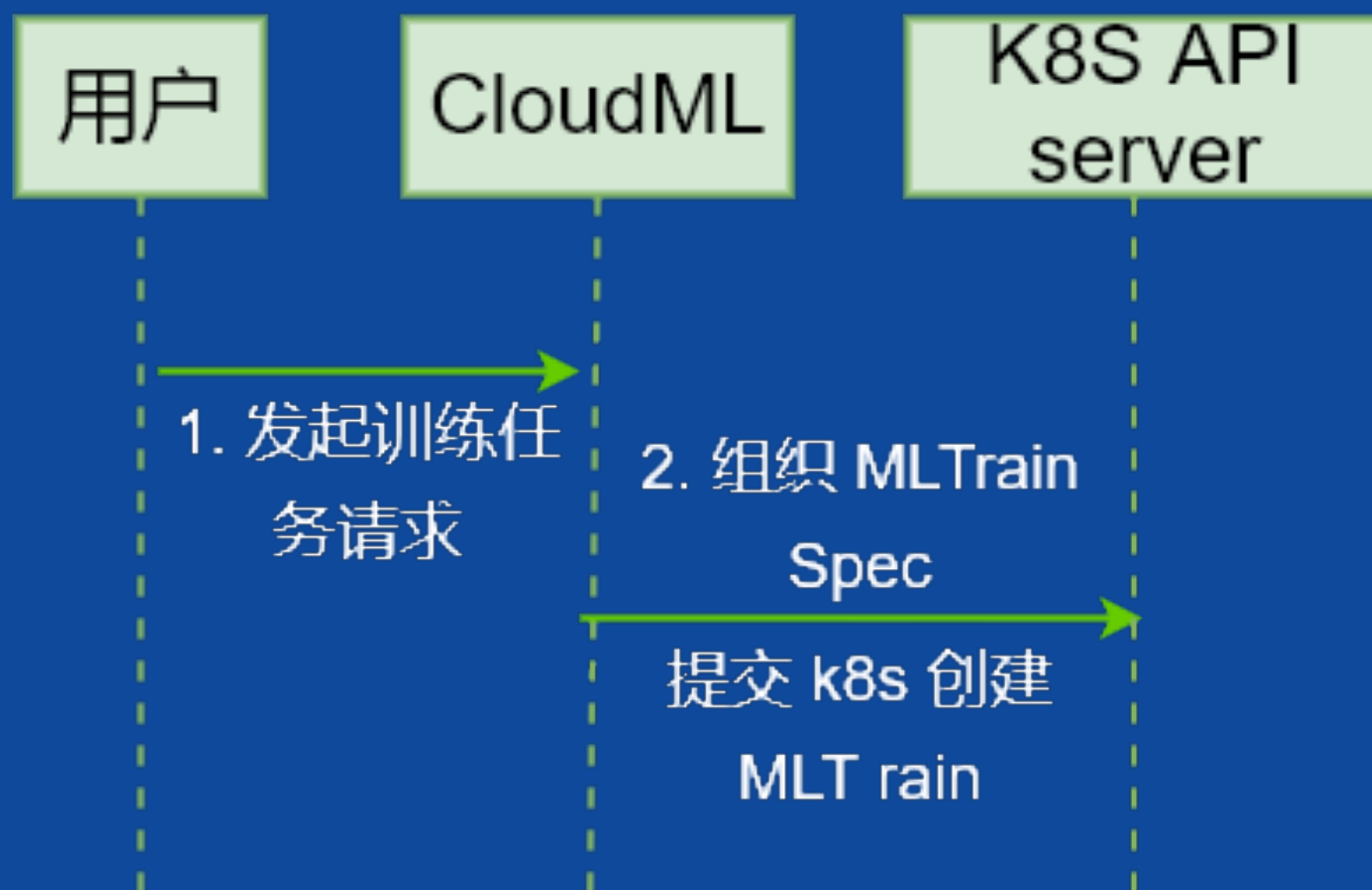
- 用户给定训练代码及启动命令（可选：定义数据及产出物）
- 支持选择不同的机器学习框架及运行时环境（CPU/GPU）
- 下发集群，需要支持分布式启动训练，且需要遵循一定的调度优化策略

# ML Engine 训练任务流程





# ML Engine 训练任务定义



```
37 // TrainSpec defines the desired state of Train
38 type TrainSpec struct {
39     // INSERT ADDITIONAL SPEC FIELDS - desired state of cluster
40     // Important: Run "make" to regenerate code after modifying this file
41     ClusterSpec map[TrainRoleName]TrainRoleSpec `json:"cluster_spec,omitempty"`
42
43     // Framework like: tf, mpi, pytorch, default: tf
44     // +optional
45     Framework string `json:"framework,omitempty"`
46     // +optional
47     Image string `json:"image,omitempty"`
48     // +optional
49     Command []string `json:"command,omitempty"`
50     // +optional
51     Args []string `json:"args,omitempty"`
52
53     // +optional
54     Env []corev1.EnvVar `json:"env,omitempty"`
55 }
56
57 // TrainRoleSpec defines the desired state of Train Cluster Role
58 // e.g. chief, worker
59 type TrainRoleSpec struct {
60     Count int32 `json:"count"`
61     // +optional
62     Image string `json:"image,omitempty"`
63     Resources corev1.ResourceRequirements `json:"resources,omitempty"`
64     // TODO(xychu): add role level cmd and args support
65     // +optional
66     Command []string `json:"command,omitempty"`
67     // +optional
68     Args []string `json:"args,omitempty"`
69
70     // +optional
71     Env []corev1.EnvVar `json:"env,omitempty"`
72
73     // NOTE(xychu): use affinity to support node selector and other placement requirements
74     // +optional
75     Affinity *corev1.Affinity `json:"affinity,omitempty"`
76 }
77 }
```

# ML Engine 训练任务定义

K8S API  
server

ML Engine  
Controller

TFJob  
Controller

MPIJob  
Controller

3. 触发 ML Engine  
Controller

4. 根据 MLTrain Spec  
决定创建 TFJob

或者是  
MPIJob

kubeflow/tf-operator

kubeflow/mapi-  
operator

...

```
// +genclient
// +k8s:deepcopy-gen:interfaces=k8s.io/apimachinery/pkg/runtime.Object
```

```
24 // +resource:path=tfjob
```

```
25
26 // Represents a TFJob resource.
```

```
27 type TFJob struct {
```

```
28     // Standard Kubernetes type
```

```
29     metav1.TypeMeta `json:",inlin
```

```
30
```

```
31     // Standard Kubernetes object
```

```
32     metav1.ObjectMeta `json:"meta
```

```
33
```

```
34     // Specification of the desired
```

```
35     Spec TFJobSpec `json:"spec,c
```

```
36
```

```
37     // Most recently observed status
```

```
38     // Read-only (modified by the
```

```
39     Status common.JobStatus `js
```

```
40 }
```

```
41
```

```
42 // TFJobSpec is a desired state description of the TFJob.
```

```
43 type TFJobSpec struct {
```

```
44     // Specifies the duration (in seconds) since startTime during which the job can
45     // before it is terminated. Must be a positive integer.
```

```
46     // This setting applies only to pods where restartPolicy is OnFailure or Always
47     // +optional
```

```
48     ActiveDeadlineSeconds *int64 `json:"activeDeadlineSeconds,omitEmpty"
```

```
49
```

```
50     // Number of retries before marking this job as failed.
```

```
51     // +optional
```

```
52     BackoffLimit *int32 `json:"backoffLimit,omitEmpty"
```

```
53
```

```
54     // Defines the policy for cleaning up pods after the TFJob completes.
```

```
55     // Defaults to Running.
```

```
56     CleanPodPolicy *common.CleanPodPolicy `json:"cleanPodPolicy,omitEmpty"
```

```
57
```

```
58     // Defines the TTL for cleaning up finished TFJobs (temporary
```

```
59     // before kubernetes adds the cleanup controller).
```

```
60     // It may take extra ReconcilePeriod seconds for the cleanup, since
```

```
61     // reconcile gets called periodically.
```

```
62     // Defaults to infinite.
```

```
63     TTLSecondsAfterFinished *int32 `json:"ttlSecondsAfterFinished,omitEmpty"
```

```
64
```

```
65     // A map of TFReplicaType (type) to ReplicaSpec (value). Specifies the TF cluster
```

```
66     // For example,
```

```
67     // {
```

```
68     //     "PS": ReplicaSpec,
```

```
69     //     "Worker": ReplicaSpec,
```

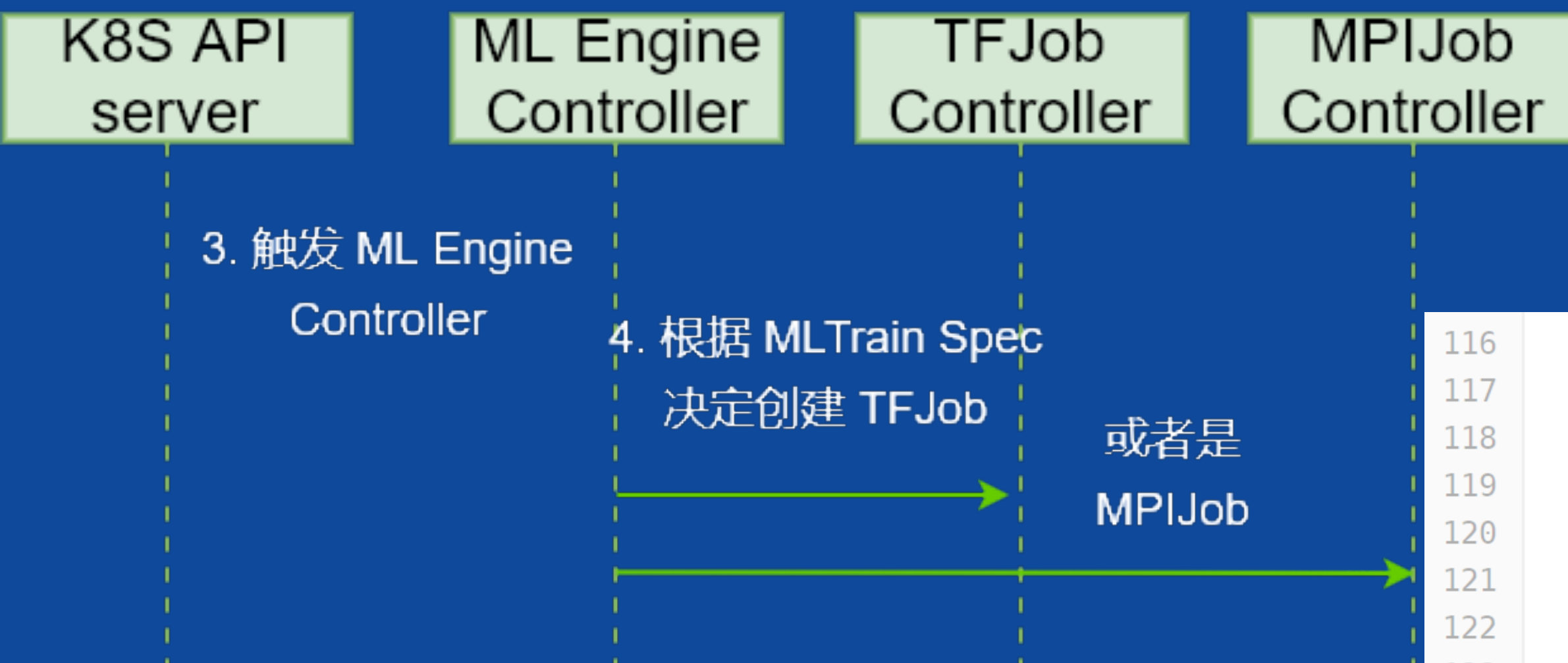
```
70     // }
```

```
71     TFReplicaSpecs map[TFReplicaType]*common.ReplicaSpec `json:"tfReplicaSpecs"
```

```
72 }
```



# ML Engine 多框架支持



## 插件化设计支持

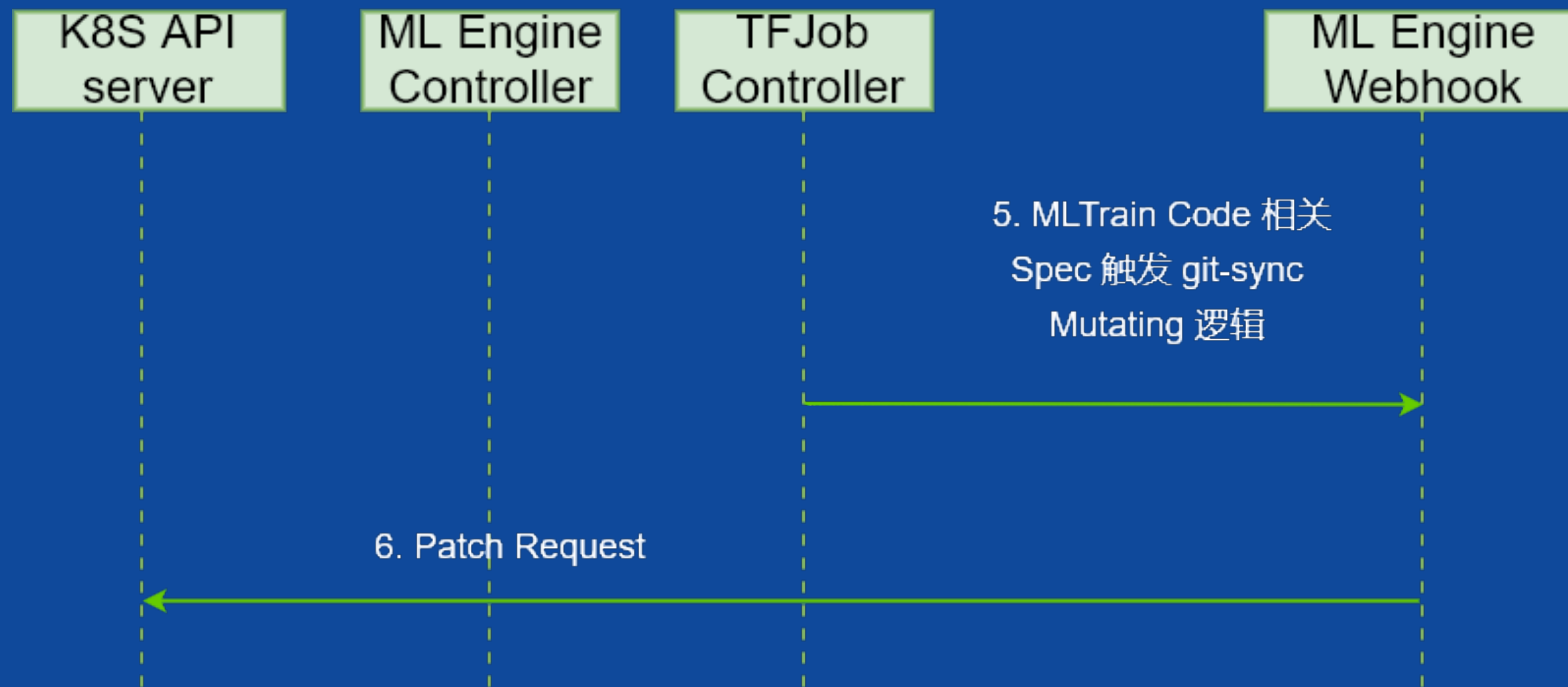
解决问题:

- 用户接口层统一
- 各个框架状态统一

```
116 // TODO(user): Change this to be the object type created by your controller
117 // Define the desired TrainJob object
118 if handler, ok := HandlerMap[instance.Spec.Framework]; ok {
119     reconcileResult, err := handler.Handle(instance, r.Client, r.scheme)
120     if err != nil {
121         return reconcileResult, err
122     }
123     // Call Update status
124     err = r.Status().Update(context.Background(), instance)
125     if err != nil {
126         return reconcile.Result{}, err
127     }
128     return reconcileResult, nil
129 } else {
130     msg := fmt.Sprintf("Unsupported Train Framework found %s\n", instance.Spec.Framework)
131     log.Info(msg, "handlers", HandlerMap)
132     err = errors.NewBadRequest(msg)
133     return reconcile.Result{}, err
134 }
```



# Train 代码拉取



## Mutating Webhook

- git-sync

<https://github.com/kubernetes/git-sync>

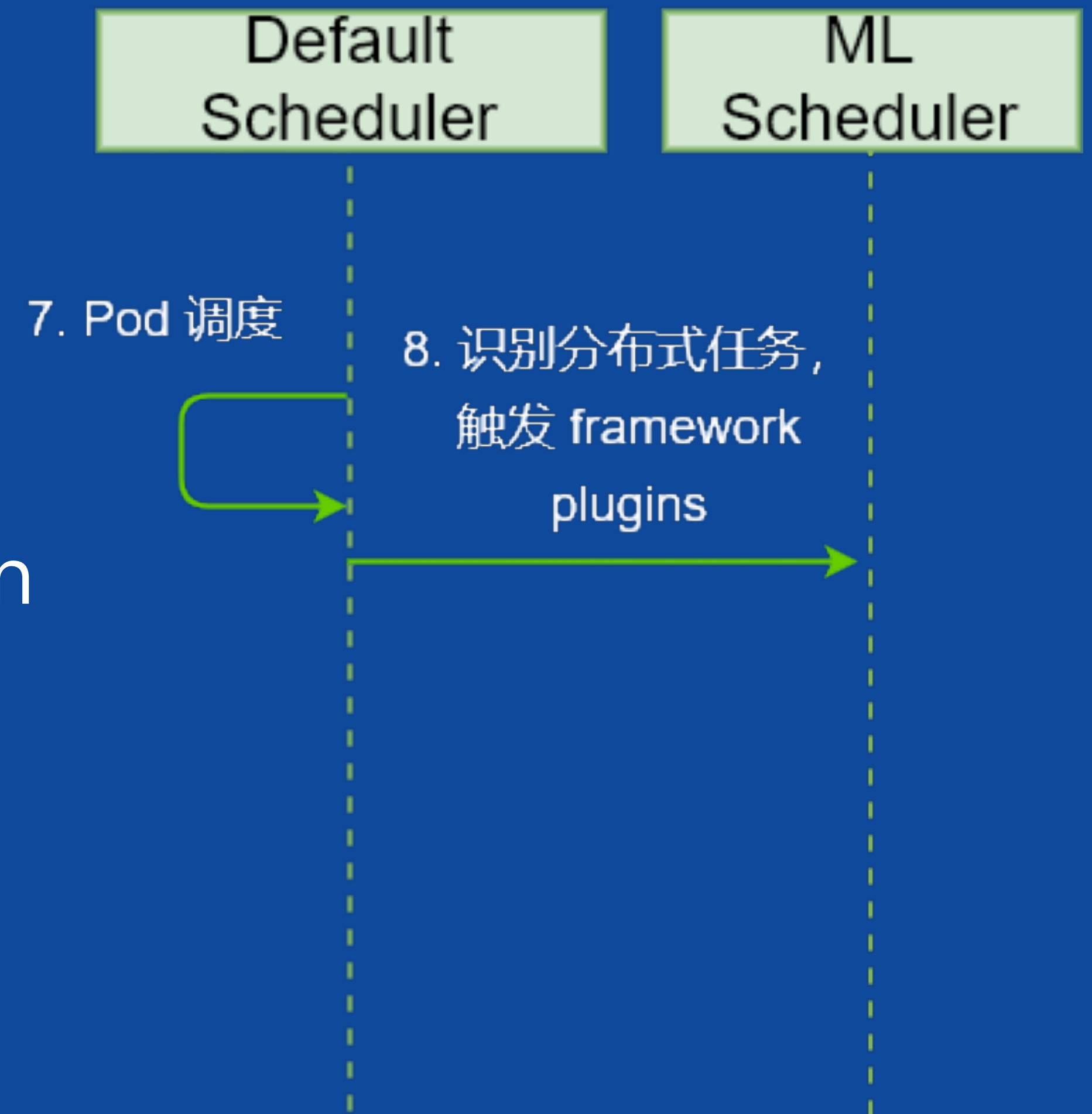
```
# ...
containers:
- name: git-sync
  image: k8s.gcr.io/git-sync:v9.3.76
  args:
    - "-ssh"
    - "-repo=git@github.com:foo/bar"
    - "-dest=bar"
    - "-branch=master"
volumeMounts:
- name: git-secret
  mountPath: /etc/git-secret
securityContext:
  runAsUser: 65533 # git-sync user
# ...
```

# ML Engine 调度优化

GPU Mosted Requested Priority

Distributed Coscheduling Framework plugin

Affinity/Anti-Affinity



# 心得总结

优点：

- K8s 原生的扩展性
- 服务的可用性保证
- 服务的易维护性

缺点：

- admission webhook 双刃剑，建议设置过滤条件
- 多框架支持，各个框架不同版本之间的 golang 依赖



# TABLE OF CONTENTS 大纲

---

- 小米 CloudML 机器学习平台简介
- ML Engine 架构设计演进
- ML Engine 对多框架的分布式训练支持详解
- 未来发展方向和具体工作

# 未来发展方向和具体工作

## 功能加强

- 考虑将计量/计费也抽象成 CRD
- 提供更方便的数据管理逻辑
- 更丰富的模型服务管理能力

时机成熟后会积极尝试开源

欢迎加入~