

cube-studio

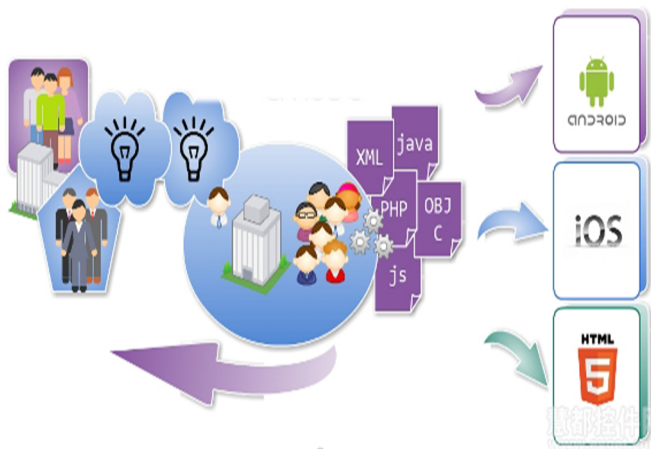
一站式云原生机器学习平台



行业现状：AI平台发展现状和痛点

现状

- ✓ 数据智能化+大模型带来AI平台刚性需求
- ✓ 智能化转型在早期，一般不具备AI平台能力
- ✓ 数据和算力都不想上云



自主可控

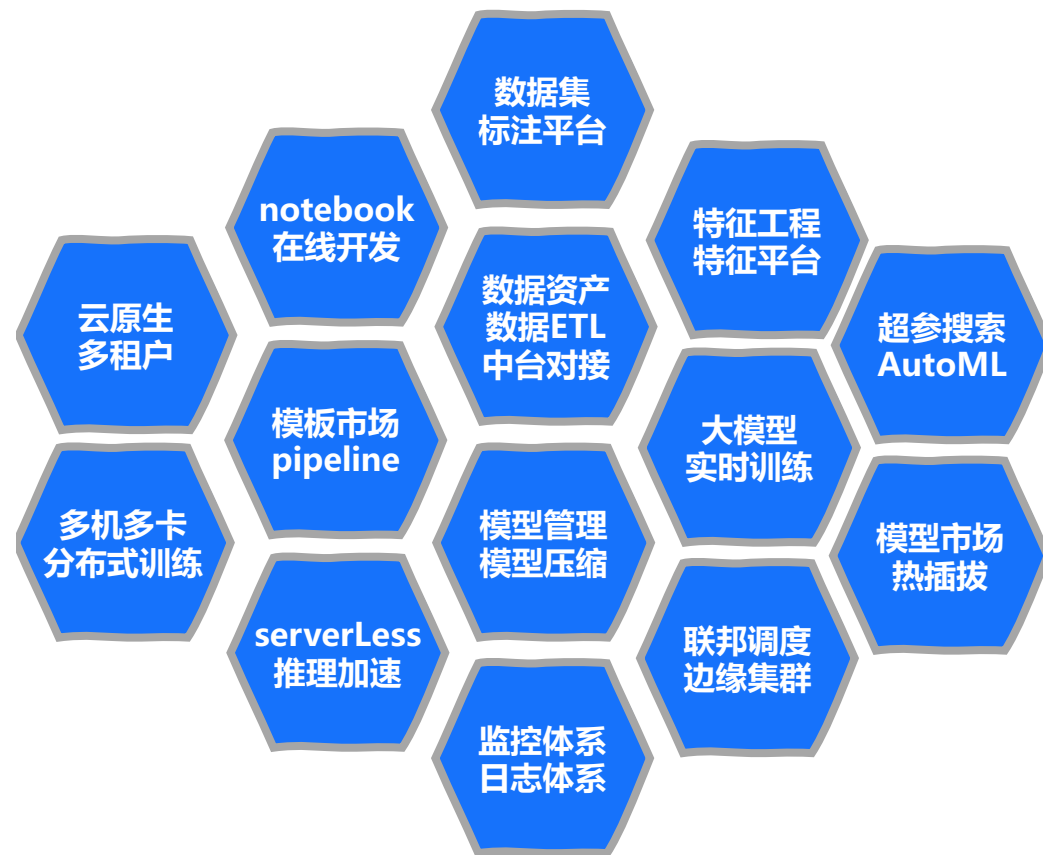
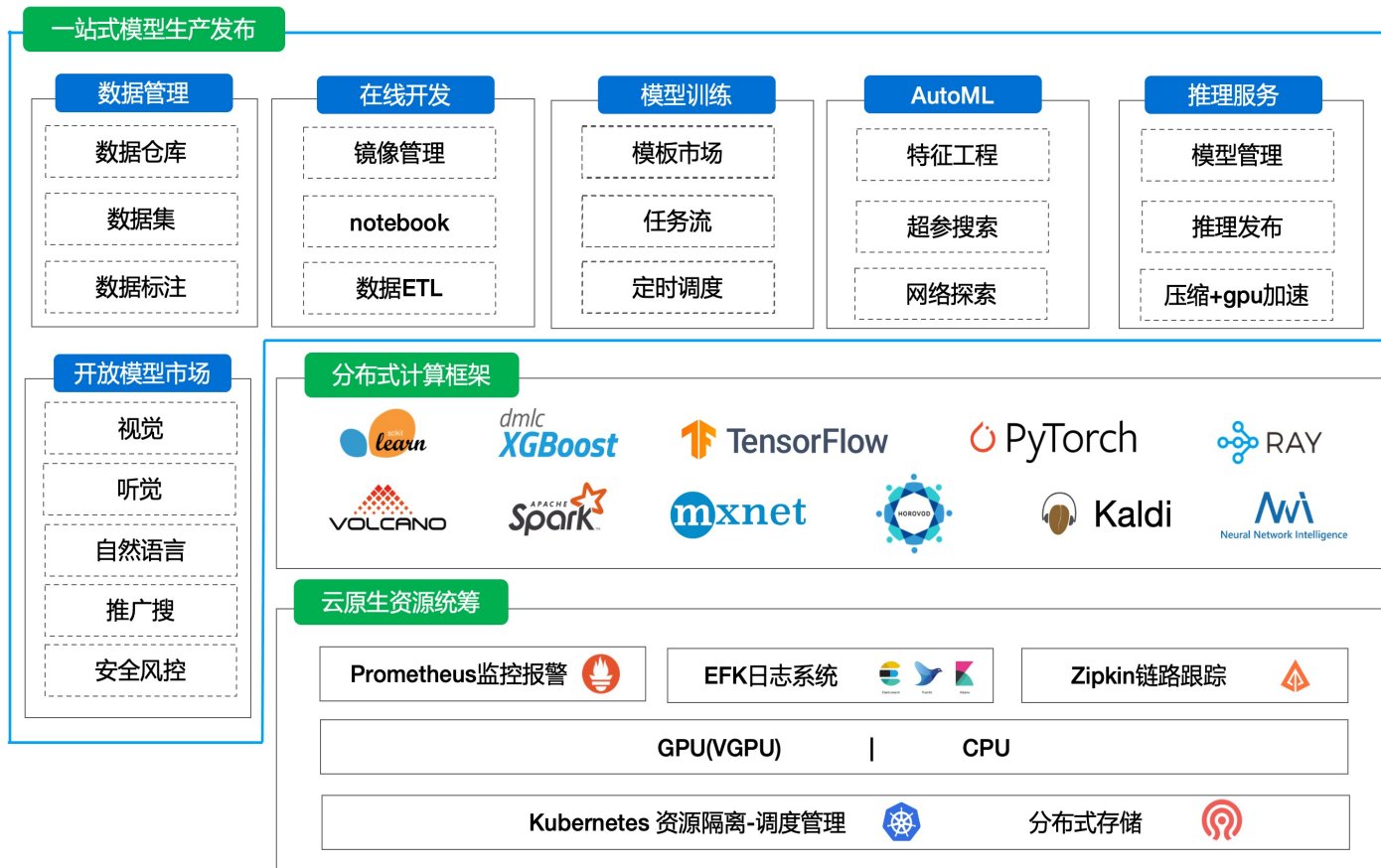
- ✓ 代码自主可控，可自主改造定制
- ✓ 私有化部署
- ✓ 算力可控，国产gpu/npu等

难点

- ✓ 专业门槛高，培养周期长
- ✓ 投入产出比不明朗
- ✓ 被破内卷，无法找到适合自己的落脚点



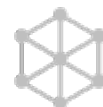
产品功能架构



AI中台一站式



数据中台对接



AIHub模型市场



gpt大模型

产品价值

90%

完备性+稳定性

快速具备完备的算法平台：cube-studio 包含完备的一站式学习平台能力，开源社区的补充，使其具备了各种场景的兼容性，腾讯内部实践，使其具备较好的稳定性。

100%

源码交付

细分市场的适配性：支持源码交付，可自行为客户进行定制，或加工改造为自己的专有产品。平台开发框架适合人力少情况的定制开发。

60%

算力需求

算力需求多样性：可私有化部署，支持国产算力，支持arm64架构，支持serverless弹性伸缩，支持边缘集群模式，支持多k8s集群。

80%

合作共赢

与解决方案商合作共赢：不直接对接细分领域客户，与解决方案商无竞争关系，协助解决方案商/算力厂商与客户项目落地，参与分成。

友商功能对标

指标	Cube studio	火山引擎	百度Paddle	kube flow	微众 Prophecis	Meta Spore	阿里 PAI
是否自主研发	✓	✓	✓	/	✓	✓	✓
功能全面性	★★★	★★★	★★★★	★★★	★★	★	★★★★
灵活性	高	中	中	低	中	中	中
兼容性	★★★★	★★★	★★★	★	★★	★★	★★★
源码交付	是	否	否	是	否	否	否
用户体验	★★★★	★★★★	★★★★	★	★★	★★	★★★★
使用成本	★	★★★	★★★★	★	★★	★★	★★★★
是否开源	开源	不开源	部分	✓	✓	✓	部分
本地化服务	★★★★	★★★	★★★	/	/	/	★★★
销售渠道	★	★★★	★★★★	/	/	/	★★★★

国内mlops github开源第一

上百家企业私有化部署

The screenshot shows the GitHub repository page for 'tencentmusic / cube-studio'. The repository is public and has 86 issues, 9 pull requests, 372 forks, and 1.4k stars. The 'About' section describes it as an open-source cloud-native one-stop machine learning/deep learning AI platform. It supports SSO login, multi-tenant/multi-project groups, data asset integration, notebook online development, task flow pipeline编排, multi-machine multi-card distributed algorithm training, hyperparameter search, inference service VGPU, multi-cluster scheduling, edge computing, serverless, labeling platform, automated labeling, data management, large model one-click fine-tuning, MLOps, private knowledge base, AI application store, one-click development/inference/fine-tuning, private deployment, and support for domestic arm64 chips. The repository includes files like aihub, docs, images, install, job-template, myapp, .gitignore, CONTRIBUTING.md, LICENSE, and README.md.



开源星「001号」落地 (cube-studio) , 开源贡献赢神秘大礼包!

滕源会 TME数据智能 2022-08-17 19:45 发表于广东

第一波入选腾讯摘星计划

今年5月6日,腾讯·滕源会社区联合(项目名)等在内的80余家开源社区、国内外开源基金会等,共同发起「开源摘星计划」。开展3月以来,我们累计为近百位优秀摘星贡献者,送出激励大礼包300余份;同时为700位的开源爱好者搭建了共同的交流乐园,帮助很多人完成了从开源萌新到「过来人」的成长、蜕变。

国内开源mlops平台第一名

日访问量3000+, 日clone量100+, 社区参与900+人

主要功能介绍-算力/存储/用户管理

算力

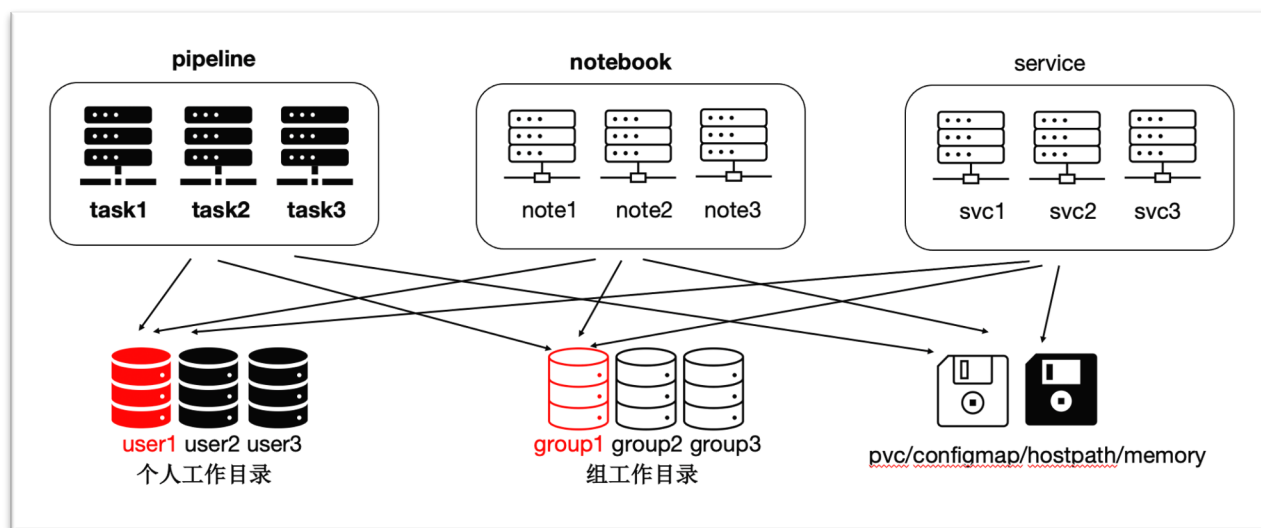
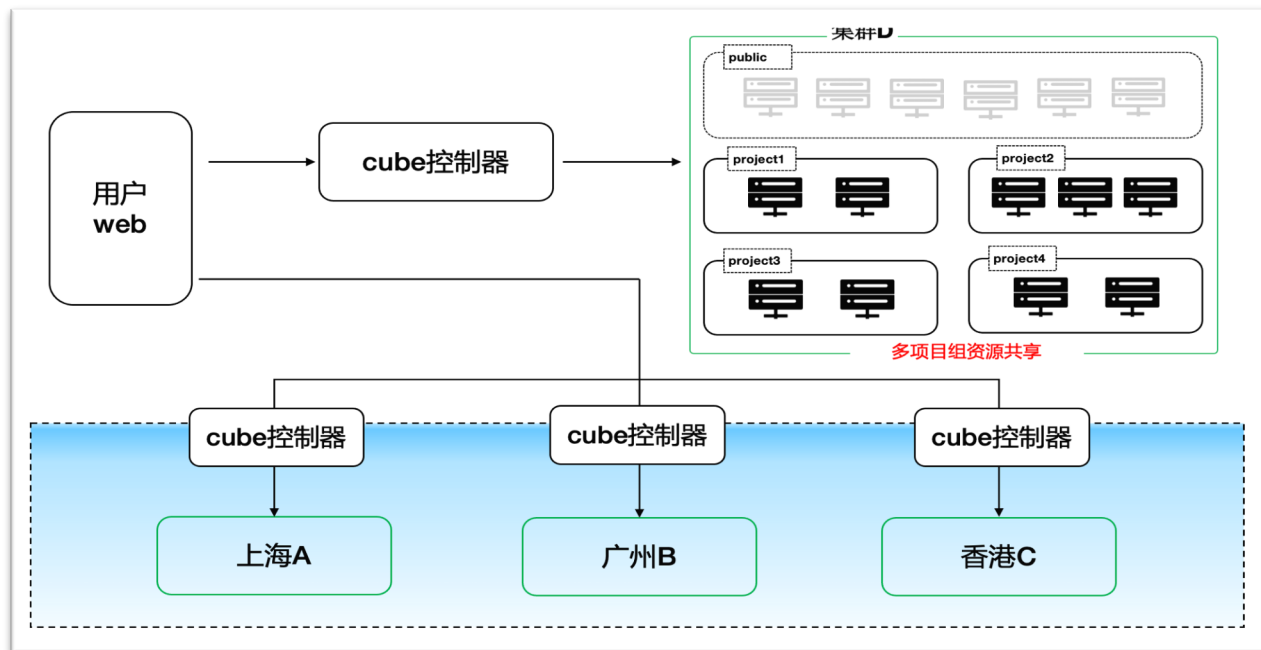
- ✓云原生统筹平台cpu/gpu等算力
- ✓支持划分**多资源组**，支持**多k8s**集群，多地部署
- ✓支持T4/V100/A100/**昇腾/VGPU**等异构GPU/**NPU**环境
- ✓支持**边缘集群**模式，支持边缘节点上开发/训练/推理
- ✓支持**鲲鹏**芯片**arm64**架构，**RDMA**

存储

- ✓自带分布式存储，支持多机分布式下文件处理
- ✓支持**外部存储挂载**，支持项目组挂载绑定
- ✓支持个人存储空间/组空间等多种形式
- ✓平台内存储空间不需要迁移

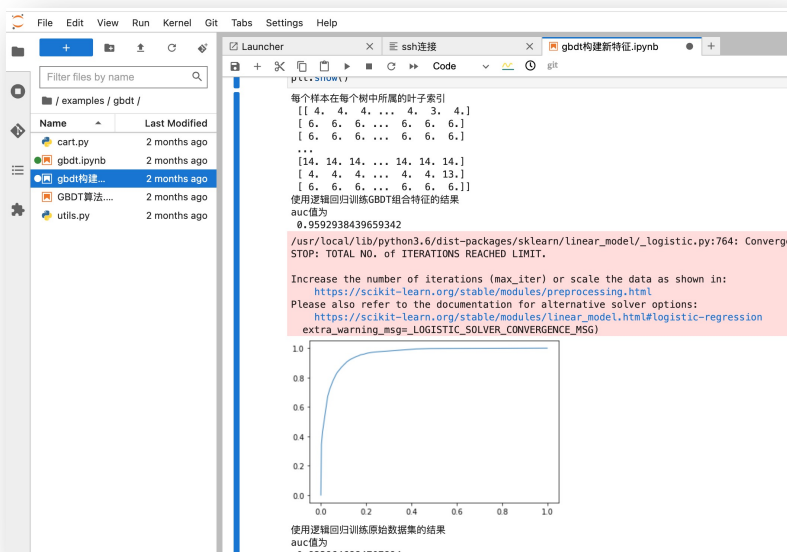
用户权限

- ✓支持**sso**登录，对接公司账号体系
- ✓支持项目组划分，支持配置相应项目组用户的权限
- ✓管理平台用户的基本信息，**组织架构**，rbac权限体系

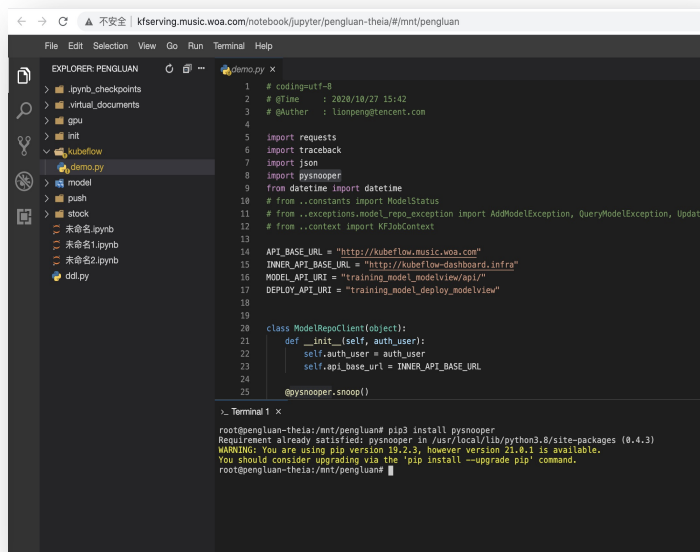


MLOPS-一站式开发工具

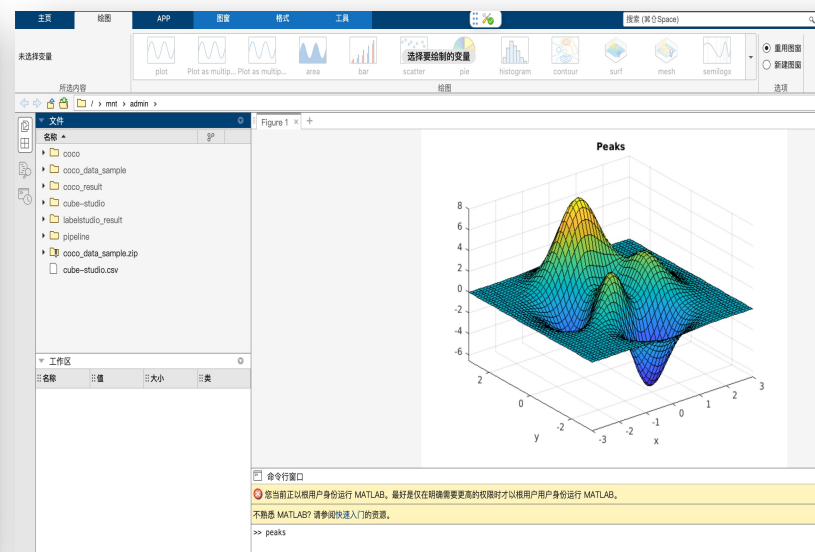
- ✓ 系统多租户/多实例管理，在线交互开发调试，无需安装三方控件，只需浏览器就能完成开发。
- ✓ 支持vscode, jupyter, Matlab, Rstudio等多种在线IDE类型
- ✓ Jupyter支持cube-studio sdk, Julia, R, python, pyspark多内核版本，



jupyter



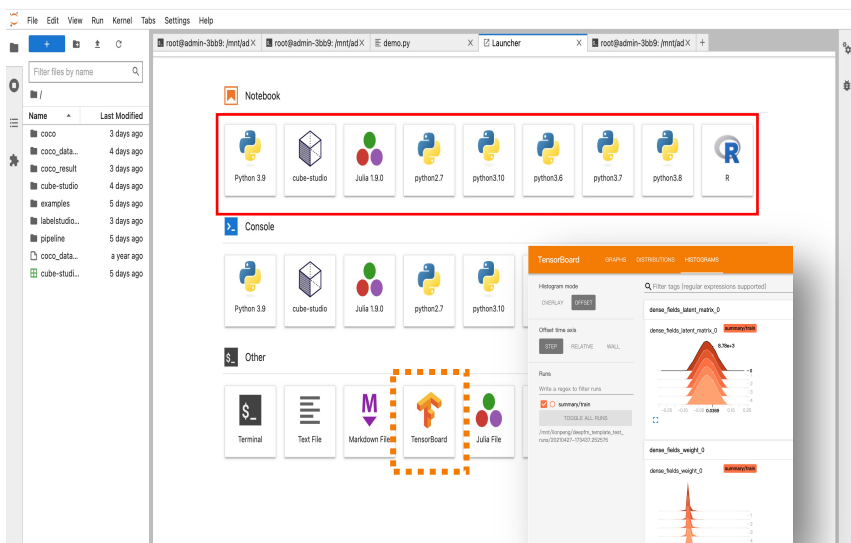
vscode



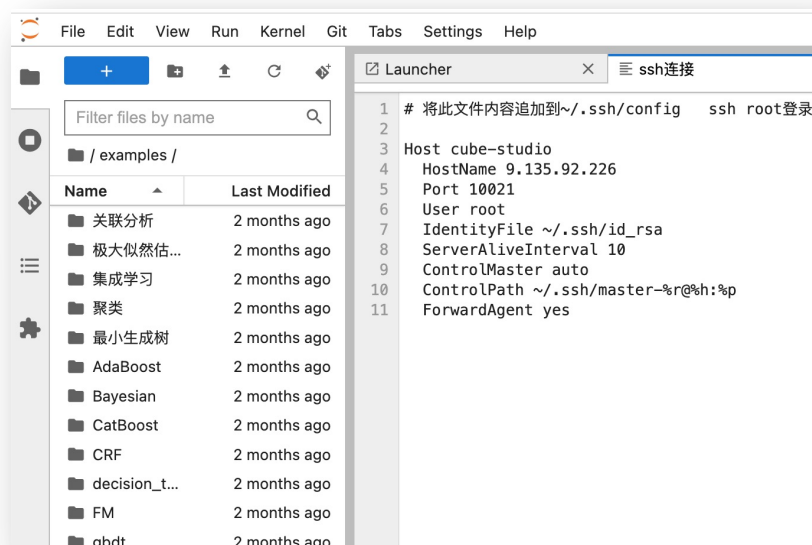
matlab

MLOPS-一站式开发工具

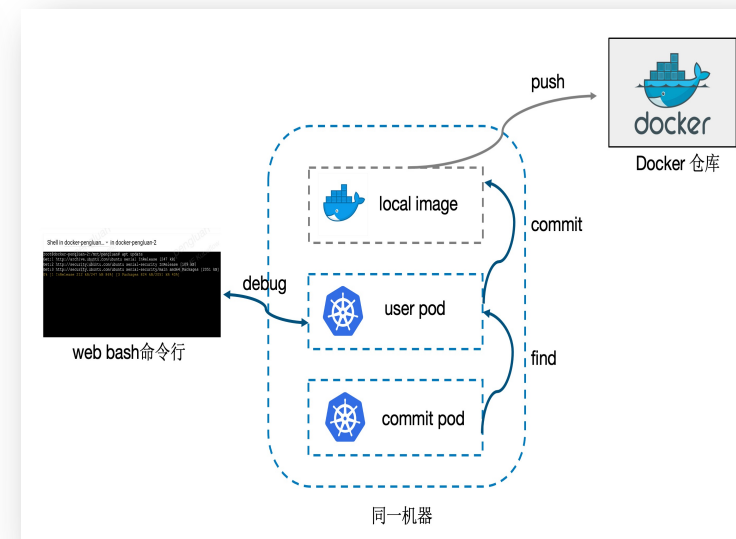
- ✓ 支持c++ , java , conda等多种开发语言 , 以及tensorboard/git/gpu监控等多种插件
- ✓ 支持ssh remote与notebook互通 , 本地进行代码开发
- ✓ 在线镜像构建 , 通过Web Shell方式在浏览器中完成构建 ; 并提供各种版本notebook , inference , gpu , python等基础镜像



多种内核/插件



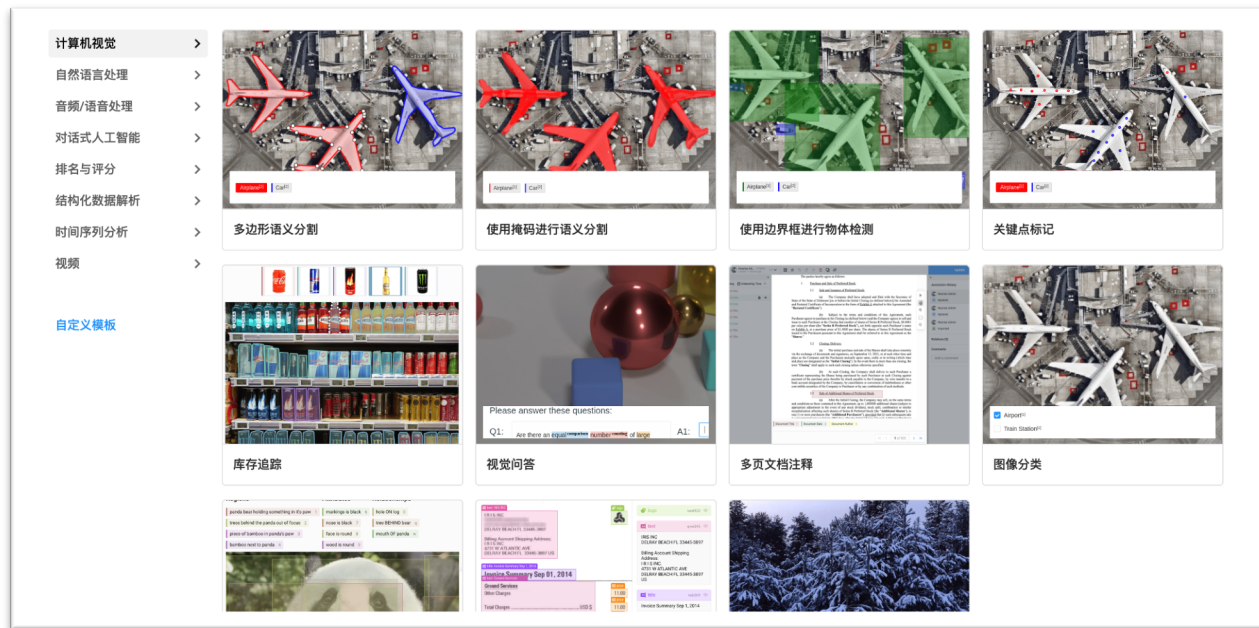
ssh远程开发



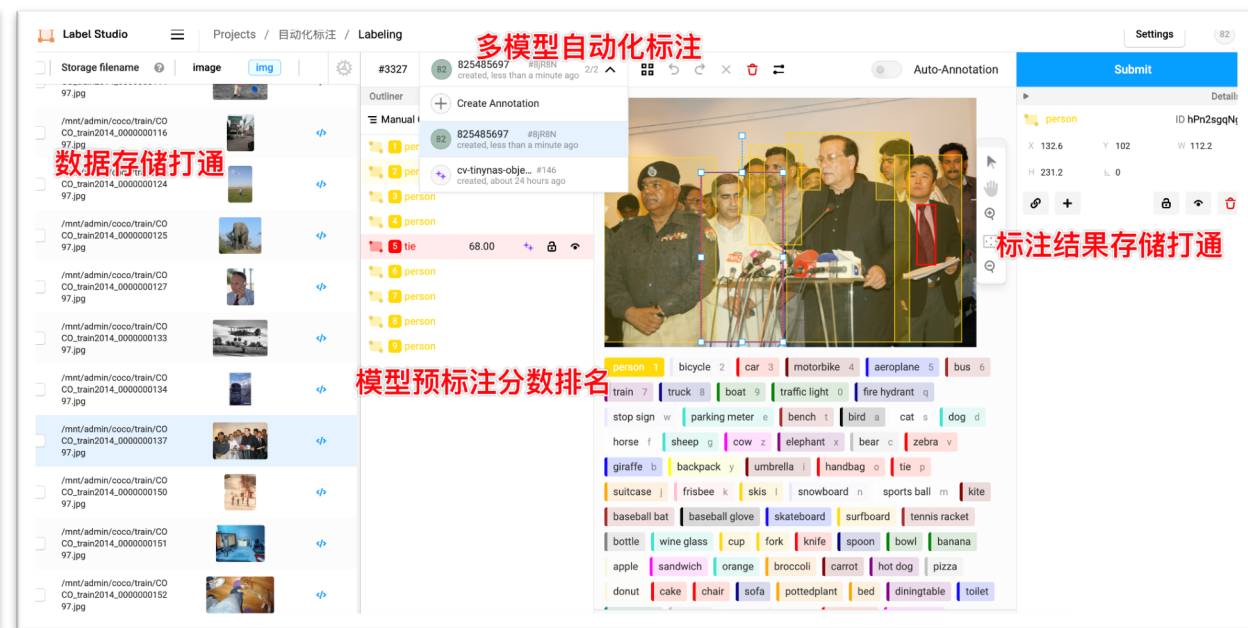
在线构建镜像

MLOPS-可视化数据标注

- ✓ 支持图/文/音/多模态/大模型多种类型标注功能，用户管理，工作任务分发
- ✓ 对接aihub模型市场，支持自动化标注；对接数据集，支持标注数据导入；对接pipeline，支持标注结果自动化训练



图文音多模态大模型类型标注



对接aihub/大模型支持自动化标注

MLOPS-拖拉拽建模流程

MI全流程

1、数据导入，数据预处理，超参搜索，模型训练，模型评估，模型压缩，模型注册，服务上线，机器学习/深度学习/大模型全流程

灵活开放

2、便捷的拖拉拽方式，编排算法dag，支持单任务和pipeline整体等多种调试运行方式
3、丰富多样的计算任务模板。支持自定义模板

分布式计算

4、便捷的多机多卡分布式训练，通过标准化方式提供分布式多机多卡计算和训练的能力，提升平台落地大规模计算所需要的效率。

模板列表

拖拉拽编排

任务配置

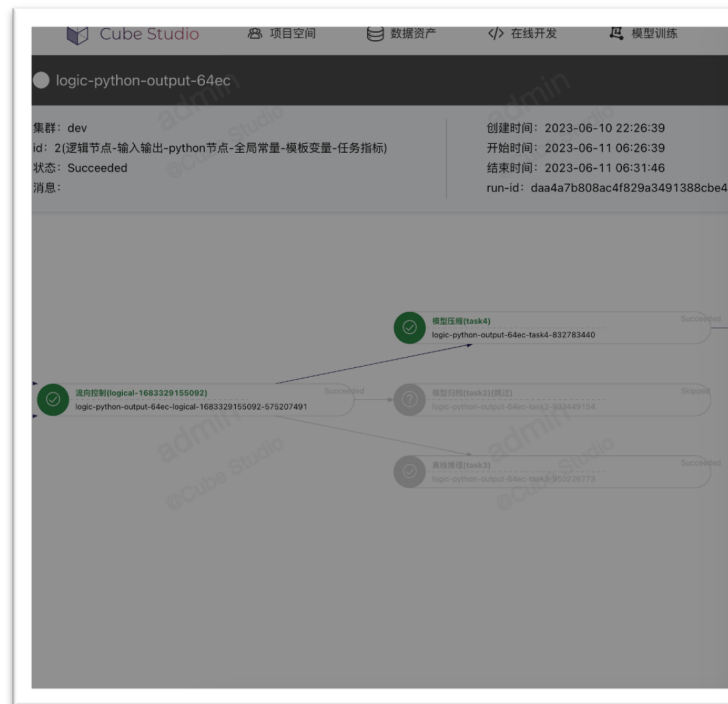
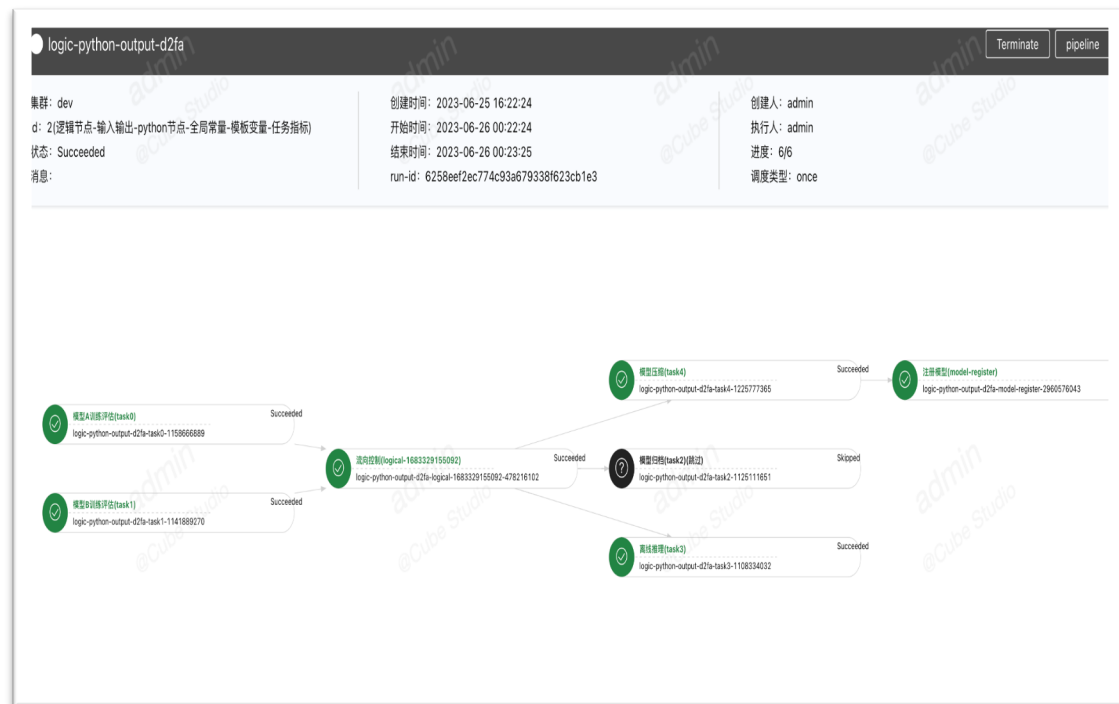
The screenshot displays the Cube Studio MLOps interface. On the left, a sidebar lists 88 assets in total, categorized into various machine learning tasks such as data import, preprocessing, training, evaluation, and deployment. The main workspace shows a drag-and-drop workflow diagram for a machine learning pipeline. The workflow includes steps like '数据导入' (Data Import), '数据处理' (Data Processing), '超参搜索' (Hyperparameter Search), '模型训练' (Model Training), '模型评估' (Model Evaluation), '模型高线推理' (Model Inference), '模型注册' (Model Registration), and '服务部署' (Service Deployment). On the right, a task configuration panel for 'model-train' is visible, showing parameters such as 'SeriousDlqin2yrs', 'save_model_dir', 'feature_columns', and 'model_params'.

MLOPS-多层次多类型算子



MLOPS-流水线调试

- ✓ Pipeline调试支持**定时**执行，支持，补录，并发限制，超时，实例依赖等。
- ✓ Pipeling运行，支持变量在任务间输入输出，全局变量，流向控制，模板变量，数据时间等
- ✓ Pipeling运行，支持**任务结果可视化**，图片、csv/json，echart源码可视化



MLOPS-超参搜索

- ✓ 界面化呈现训练各组数据，通过图形界面进行直观呈现。
- ✓ 减少以往开发调参过程的枯燥感，让整个调参过程更加生动具有趣味性，完全无需丰富经验就能实现更精准的参数控制调节。

上报当前迭代目标值

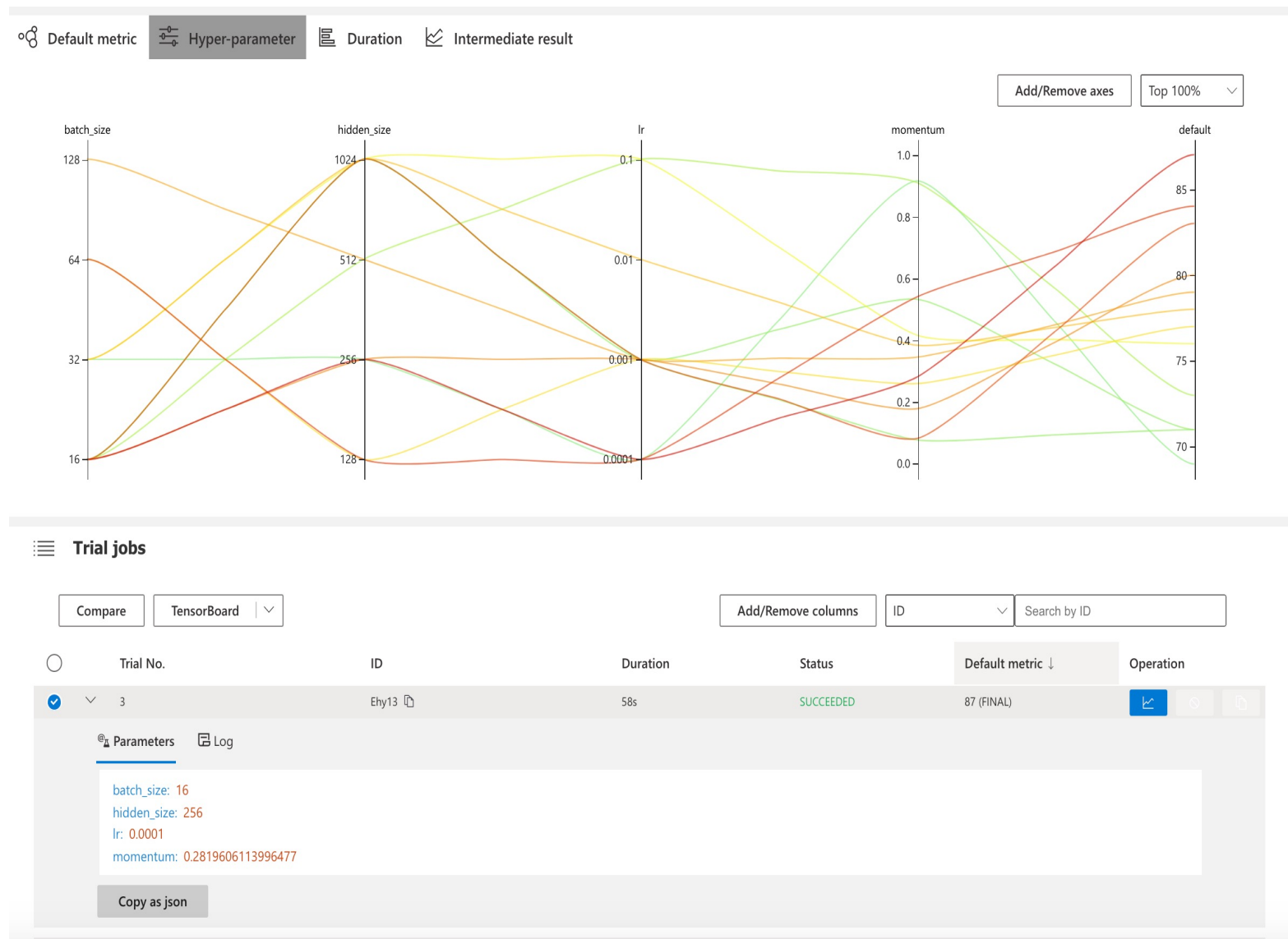
```
nni.report_intermediate_result(test_acc)
```

上报最终目标值

```
nni.report_final_result(test_acc)
```

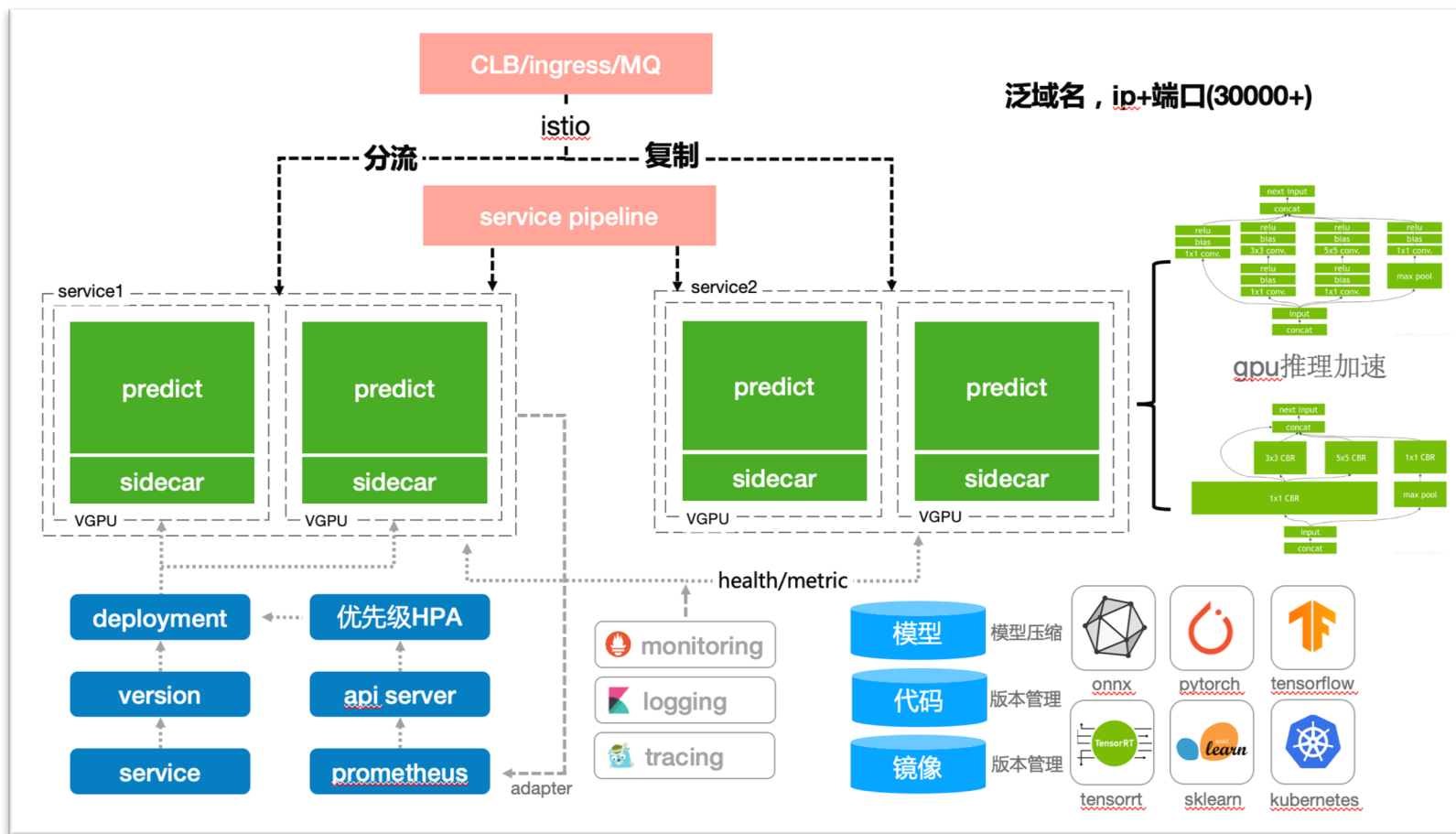
接收超参数为输入参数

```
parser.add_argument('--batch_size', type=int)
```



MLOPS-自动化零代码推理部署

- ✓ 支持模型管理注册，灰度发布，版本回退，模型指标可视化，以及在pipeline中进行模型注册
- ✓ 推理服务支持多集群，多资源组，异构gpu环境，平台资源统筹监控，**VGPU**，服务流量分流，复制，sidecar
- ✓ 支持0代码的模型发布，gpu推理加速，支持训练推理混部，服务优先级，自定义指标弹性伸缩。



MLOPS-监控报警体系

replication Controllers

Stateful Sets

服务 **N**

Ingresses

Services

配置和存储

Config Maps **N**

Persistent Volume Claims **N**

Secrets **N**

Storage Classes

集群

Pods

名称	标签	节点	状态	重启	CPU 使用率 (cores)	内存使用 (bytes)	创建时间 ↑
✓ tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95-worker-7	app: tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95 显示所有	node009138244086	Completed	0	997.00m	8.08Gi	10 hours ago
✓ tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95-worker-6	app: tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95 显示所有	node009138244122	Completed	0	692.00m	9.08Gi	10 hours ago
✓ tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95-worker-4	app: tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95 显示所有	node009138247186	Completed	0	1.00	6.68Gi	10 hours ago
✓ tfjob-6eb376c6-8b78-11eb-9420-h6rd4fe968d95-worker-1	app: tfjob-6eb376c6-8b78-11eb-9420-b6d4fe968d95 显示所有	node009138244086	Completed	0	1.00	8.78Gi	10 hours ago



运行通知

pipeline: crontab-standalone-train(调度测试)
 namespace: pipeline
 status: Running
 start_time: 2021-05-30 03:01:02
 finish_time:
[pod详情](#)

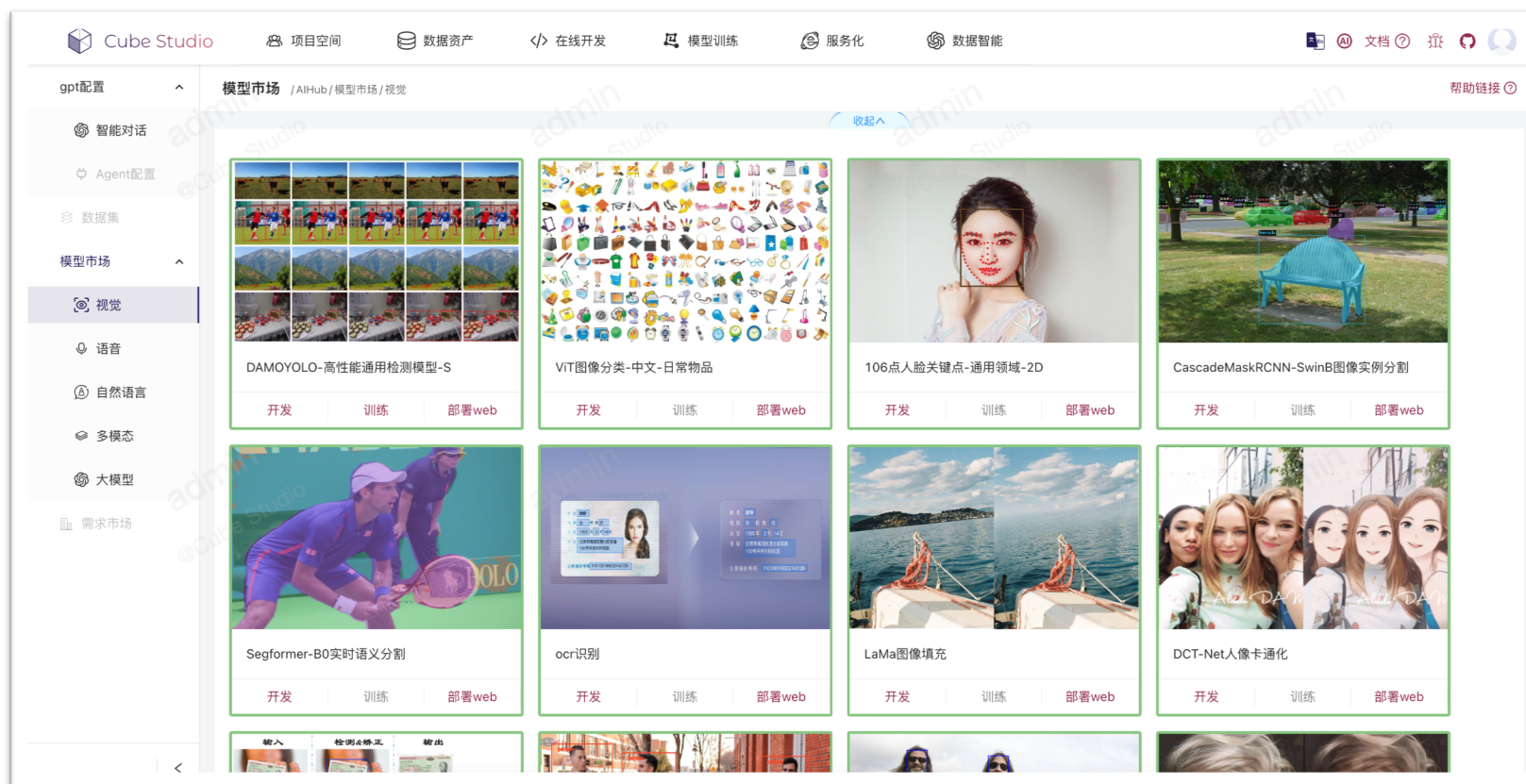
星期日 上午 4:01

workflow: crontab-standalone-train-bgm2x
 pipeline: crontab-standalone-train(调度测试)
 namespace: pipeline
 status: Succeeded
 start_time: 2021-05-30 03:01:02
 finish_time: 2021-05-30 04:01:11
[pod详情](#)

调度测试, 各task耗时, 酌情优化:
 tlinux 命令测试:1.0(h)
 tlinux 命令测试 2:1.0(h)

AIHub-模型市场

- ✓ 系统自带通用模型数量**400+**，覆盖绝大多数行业场景，根据需求可以不断扩充。
- ✓ 模型开源、按需定制，方便快捷集成，满足用户业务增长及二次开发升级。
- ✓ 模型标准化开发管理，大幅降低使用门槛，开发周期时长平均下降30%以上。



AIHub-模型市场

- ✓ AIHub模型可**一键部署**为WEB端应用，手机端/PC端皆可，实时查看模型应用效果
- ✓ 点击模型开发即可进入notebook进行模型代码的二次开发，实现**一键开发**
- ✓ 点击训练即可加入自己的数据进行**一键微调**，使模型更贴合自身场景



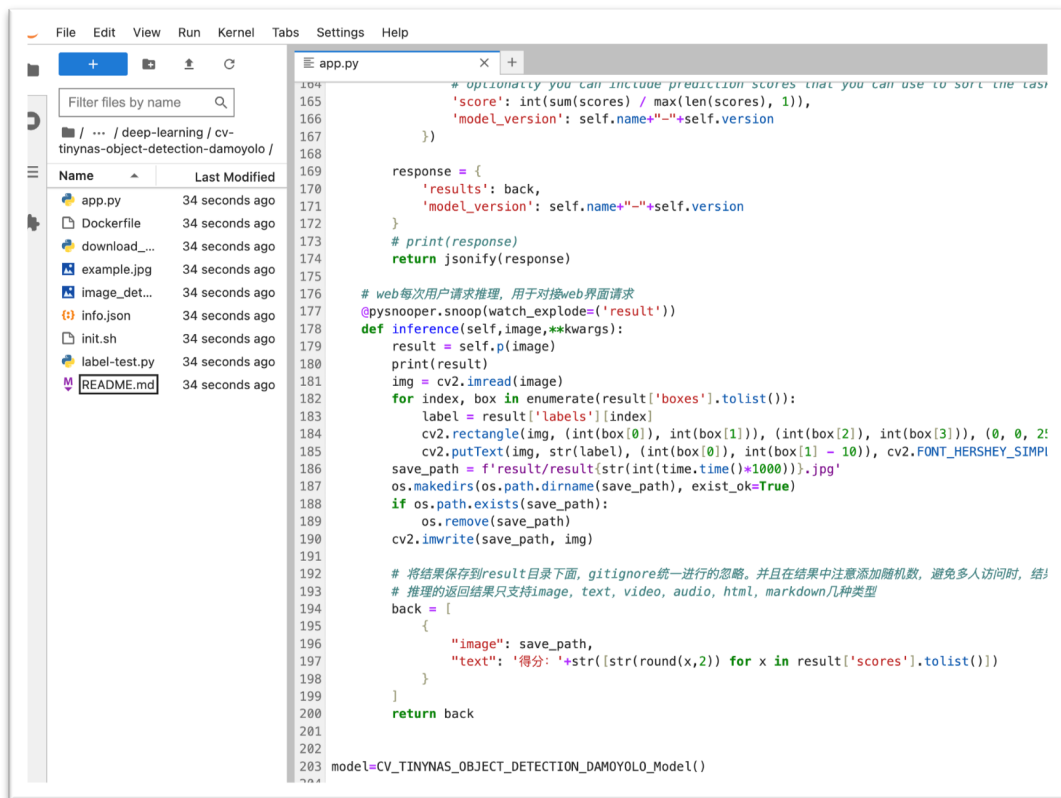
PC版



手机/微信版

AIHub-模型市场

- ✓ AIHub模型可**一键部署**为WEB端应用，手机端/PC端皆可，实时查看模型应用效果
- ✓ 点击模型开发即可进入notebook进行模型代码的二次开发，实现**一键开发**
- ✓ 点击训练即可加入自己的数据进行**一键微调**，使模型更贴合自身场景



```
File Edit View Run Kernel Tabs Settings Help
+
Filter files by name
/ ... / deep-learning / cv-tinytas-object-detection-damoyolo /
Name Last Modified
app.py 34 seconds ago
Dockerfile 34 seconds ago
download_... 34 seconds ago
example.jpg 34 seconds ago
image_det... 34 seconds ago
info.json 34 seconds ago
init.sh 34 seconds ago
label-test.py 34 seconds ago
README.md 34 seconds ago

app.py
104 # optionally you can include prediction scores that you can use to sort the task
105
106     'score': int(sum(scores) / max(len(scores), 1)),
107     'model_version': self.name+"-"+self.version
108 )}
109
110 response = {
111     'results': back,
112     'model_version': self.name+"-"+self.version
113 }
114
115 # print(response)
116
117 return jsonify(response)
118
119 # web 每次用户请求推理, 用于对接web界面请求
120 @pysnooper.snoop(watch_explode=('result'))
121 def inference(self, image, **kwargs):
122     result = self.p(image)
123     print(result)
124     img = cv2.imread(image)
125     for index, box in enumerate(result['boxes'].tolist()):
126         label = result['labels'][index]
127         cv2.rectangle(img, (int(box[0]), int(box[1]), int(box[2]), int(box[3])), (0, 0, 2:
128         cv2.putText(img, str(label), (int(box[0]), int(box[1] - 10)), cv2.FONT_HERSHEY_SIMP
129         save_path = f'result/result/{str(int(time.time()*1000)).jpg}'
130         os.makedirs(os.path.dirname(save_path), exist_ok=True)
131         if os.path.exists(save_path):
132             os.remove(save_path)
133         cv2.imwrite(save_path, img)
134
135 # 将结果保存到result目录下, gitignore统一进行的忽略, 并且在结果中注意添加随机数, 避免多人访问时, 结
136 # 推理的返回结果只支持image, text, video, audio, html, markdown几种类型
137 back = {
138     {
139         "image": save_path,
140         "text": '得分: '+str([str(round(x,2)) for x in result['scores'].tolist()])
141     }
142 }
143
144 return back
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2
```

GPT-自动化标注

✓ 借助大模型进行常规nlp任务的自动化标注



GPT-大模型微调

- ✓ 支持tfjob/pytorch/paddlejob/mindspore/mxnet等分布式训练框架
- ✓ 支持deepspeed/megatron/colossalai/horovod/mpi等分布式加速框架

The screenshot displays the Cube Studio interface for configuring a fine-tuning job. The main workspace shows a workflow with three steps: 1. 微调chatglm3 (Fine-tune chatglm3), 2. 合并lora模型 (Merge lora model), and 3. 部署chatglm3 (Deploy chatglm3). The right-hand panel is titled 'chatglm3-lora' and contains the following configuration fields:

- 内存申请 ***: 20G
- CPU申请 ***: 10
- GPU申请**: 1(A100)
- RDMA申请 ***: 1
- 超时中断**: 0
- 重试次数**: (empty)

On the left, a sidebar lists various assets under '深度学习' (Deep Learning) and '分布式加速' (Distributed Acceleration). The '深度学习' section includes tfjob, pytorchjob, paddlejob, mindspore, and mxnet. The '分布式加速' section includes mpi, colossalai, deepspeed, megatron, and horovod. A '大模型' (Large Model) section is also visible, listing llama2, chatglm3, chatglm2, and baichuan2.

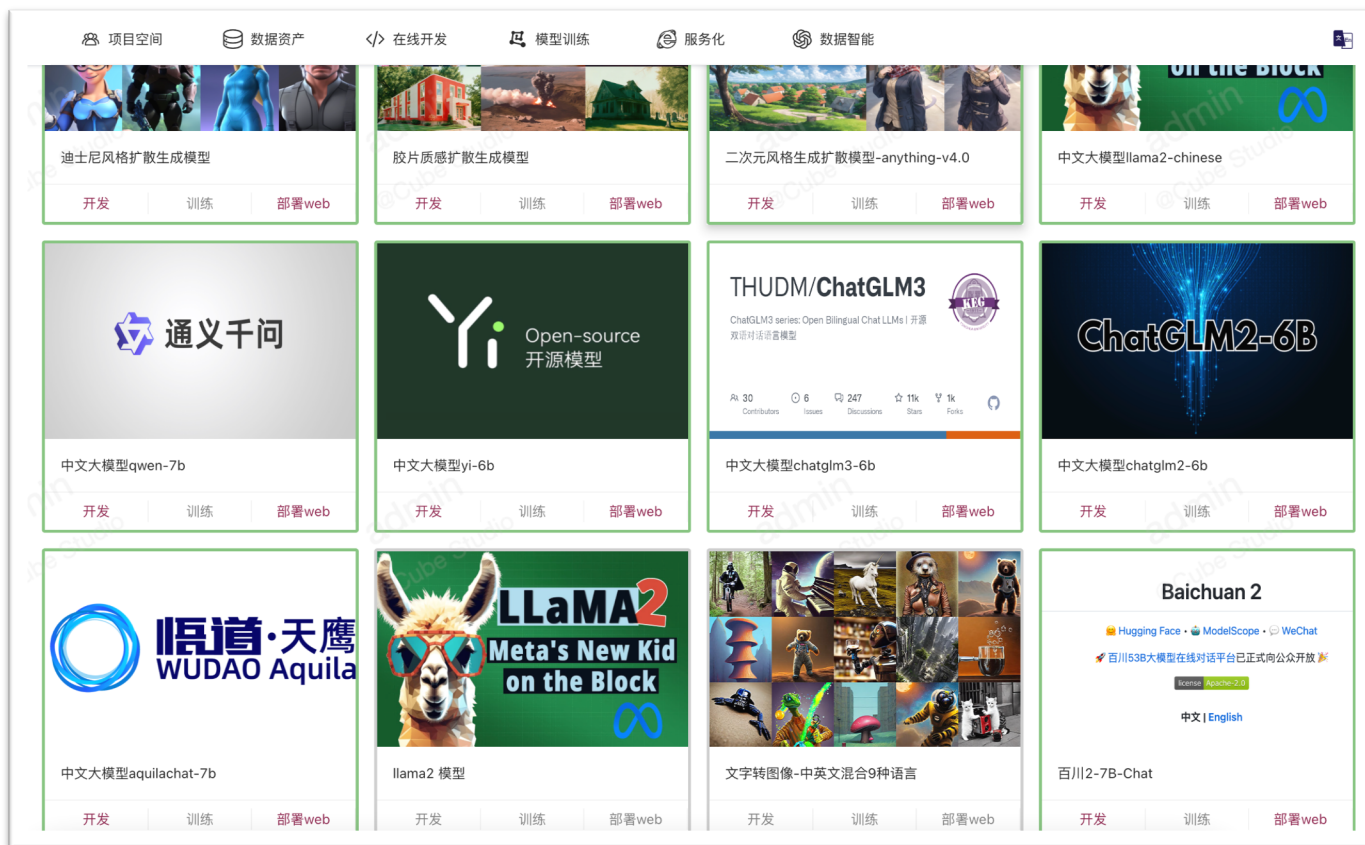
GPT-大模型微调

- ✓ 支持chatglm2/chatglm3/baichuan2/llama2等lora微调，提供chatglm微调/模型合并/部署全链路
- ✓ 支持RDMA，支持ib交换机，国产gpu卡，不同gpu卡型选择

The screenshot displays the Cube Studio interface for configuring a chatglm3-lora workflow. The main workspace shows a vertical flow of three steps: '微调chatglm3' (Fine-tune chatglm3), '合并lora模型' (Merge lora model), and '部署chatglm3' (Deploy chatglm3). The '微调chatglm3' step is currently selected, and its configuration panel on the right is visible. This panel includes fields for '内存申请' (20G), 'CPU申请' (10), 'GPU申请' (1(A100)), and 'RDMA申请' (1). The 'GPU申请' field is highlighted with a red box, indicating the selection of A100 GPUs. The left sidebar shows a list of assets, with the '大模型' (Large Model) category expanded to show options like llama2, chatglm3, chatglm2, and baichuan2. The 'chatglm3' option is also highlighted with a red box.

GPT-推理部署

- ✓ AIHub形式部署开源大模型的部署
- ✓ 支持fastchat+vllm推理加速，形成openai接口



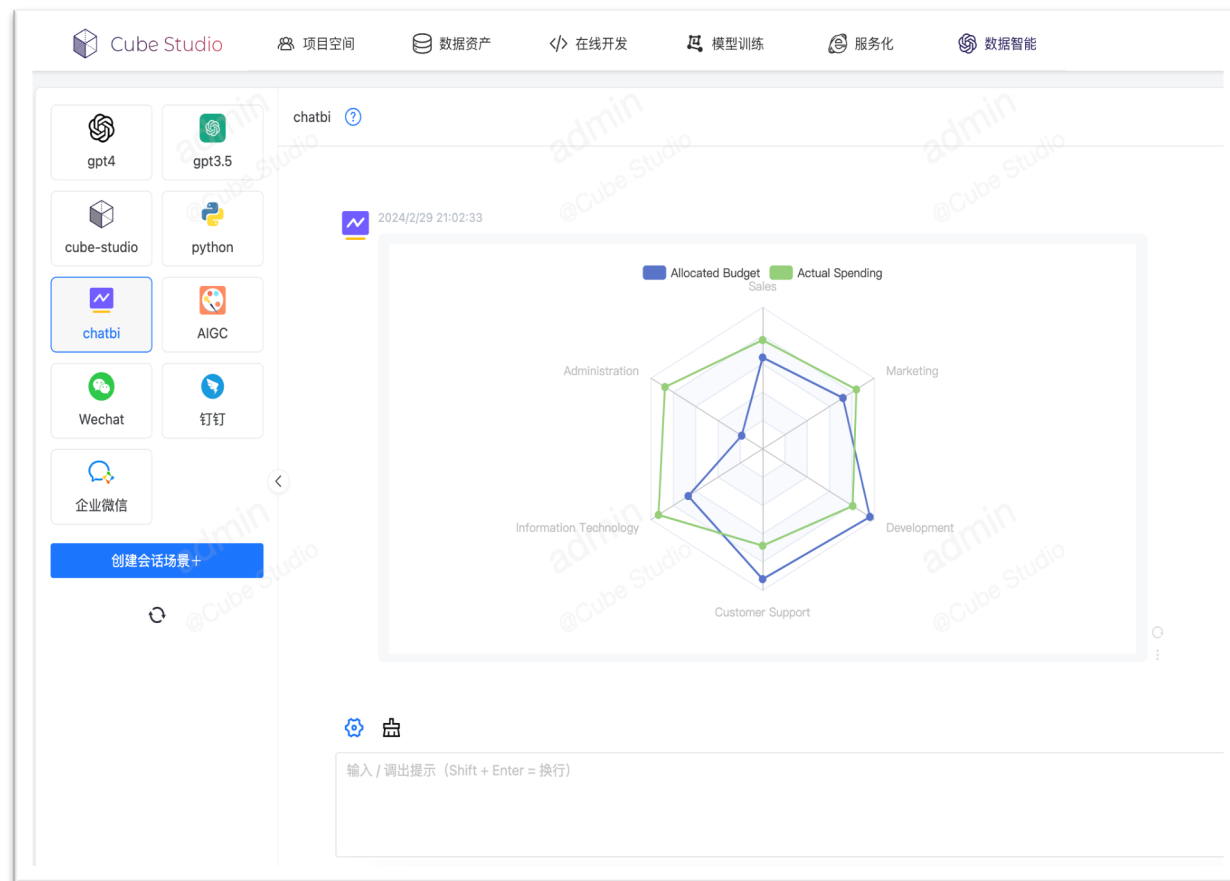
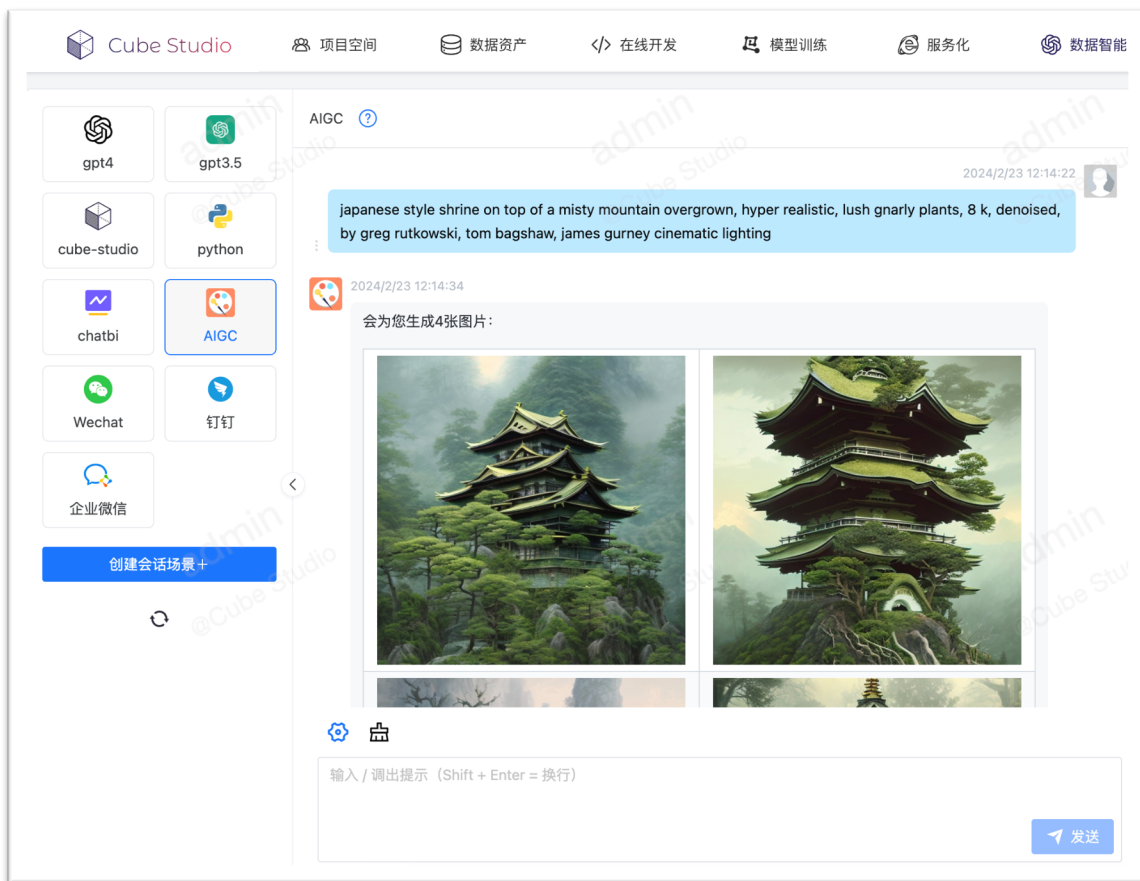
GPT-私有知识库

- ✓ 数据智能模块可配置专业领域智能对话，快速敏捷使用llm
- ✓ 可为某个聊天场景配置私有知识库文件，支持主题分割，语义embedding，意图识别，概要提取，多路召回，排序，多种功能融合



GPT-AIGC与chatbi

- ✓ AIGC：可以将智能会话与AIHub相结合，例如下面AIGC模型与聊天会话
- ✓ chatBI：智能问询，看板形式展示数据结果
- ✓ 机器人：智能会话与微信公众号/企业微信/钉钉机器人打通，可在机器人中使用知识库



数据中台对接

为了加速AI算法平台的使用，cube-studio支持对接公司原有数据中台，包括数据计算引擎sqllab，元数据管理，指标管理，维表管理，数据ETL，数据集管理

The screenshot shows the Cube Studio SQL editor interface. At the top, there are tabs for '新查询 1' and '新查询 2'. The main area contains a SQL query:

```
1 select
2 substr(imp_date,1,8) dt,
3 split(split(split(unbase64(content),'>')[0],''fontName:''[1],''')[0],
4 count(1) pv,
5 count(distinct(u)) uv
6 from
7 database: log4_hi
8 where substr(imp_date,1,8) = 20230226
9 and plat = 'ip' and regexp_replace(ver, '\\.', '') >= 10430
```

Below the query, there are controls for '正常模式', '集群: 集团', '应用组: mysql+pymysql://root:admin@mysql-service.infra:3306/kq', and a '运行' button. The '结果' section shows a progress bar for '子任务17' with steps: '准备开始', '解析', '执行', and '输出结果'. A table below shows the execution details:

子任务	开始时间	运行时长	状态	操作
select * from	2023-02-16 20:02:27	10.77秒	success	详情 结果 下载

数据资产开发工具

The screenshot shows the Cube Studio ETL workflow visualization interface. The main area displays a flowchart with nodes: '数据导入' (Data Import), '数据入库' (Data Ingestion), '局部特征计算' (Local Feature Calculation), '结果计算' (Result Calculation), and '数据导出' (Data Export). The flow starts with '数据导入' leading to '数据入库', which then branches into three parallel '局部特征计算' nodes (SQL, SparkScala, pyspark). These nodes converge into '结果计算', which finally leads to '数据导出'. The right sidebar contains configuration options for '别名', '任务元数据', '自依赖判断', and '队列'. The bottom right corner shows '监控配置' with '报警用户' set to 'admin'.

数据ETL可视化链路

三种部署模式灵活选择

针对企业和家庭的用户需求，根据不同场景对**计算实时性**的不同需求，可以提供三种建设模式

- **模式一：私有化部署**——对数据安全要求高、预算充足、自己有开发能力
- **模式二：边缘集群部署**——算力分散，多个子网环境的场景，或边缘设备场景
- **模式三：serverless集群**——成本有限，按需申请算力的场景

实时性要求低或
管理性需求场景
(公有云服务)



实时性要求较高场景
(私有化部署)

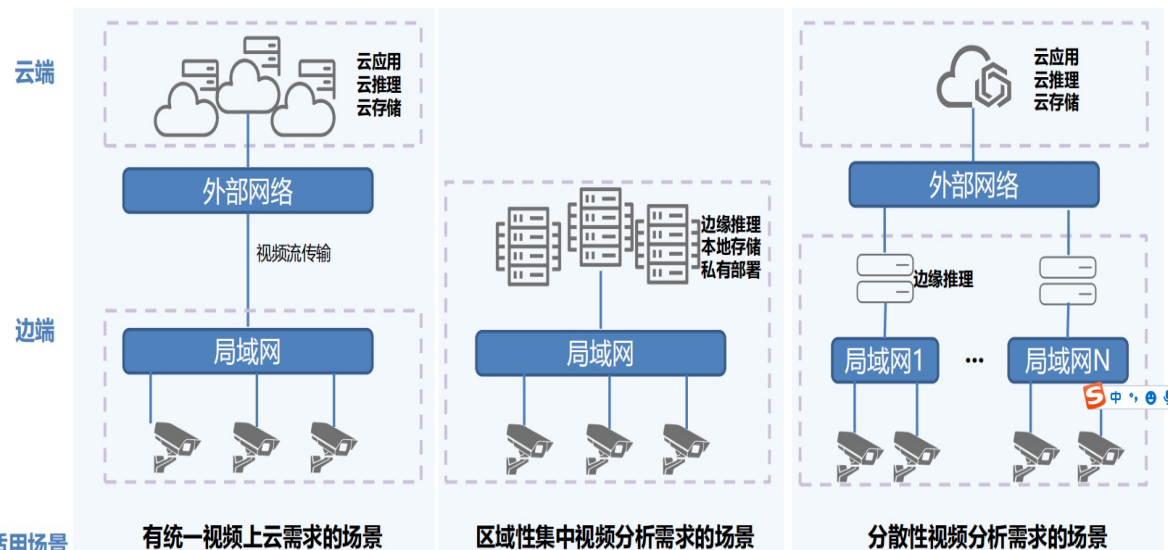


准实时流场景
(混合部署)



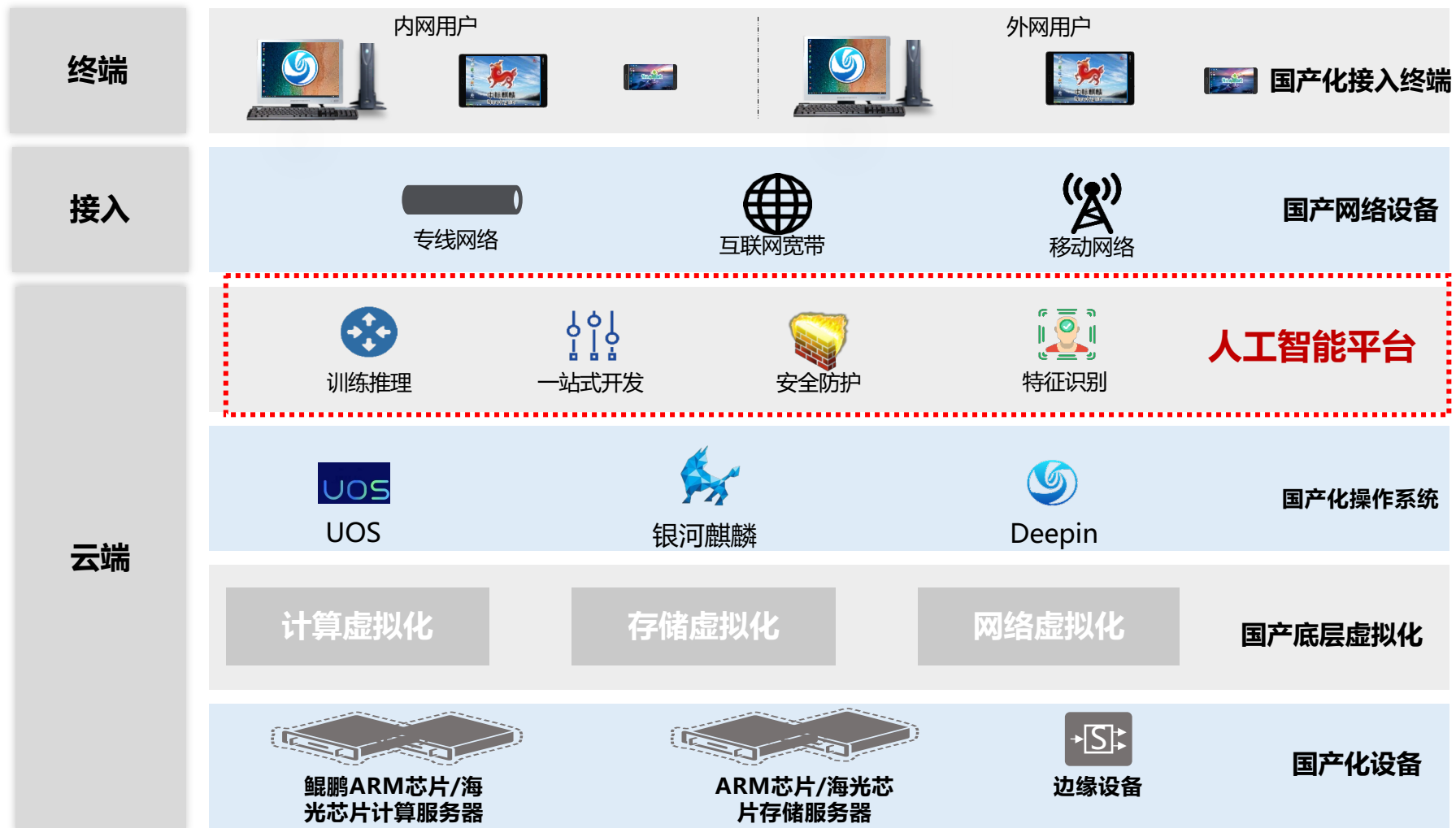
三种模式对比

类型	模式一	模式二	模式三
计算实时性	低	中	高
部署难度	低	中	高
成本	低	中	低



类型	模式一	模式二	模式三
抽帧解码能力平台	不需要	需要	需要
AI计算资源	无需	需要	需要
带宽	需要	需要	需要

国产化部署支持



产品在国产化链路中的定位

平台原生支持国产化平台部署，以满足客户对国产化产品安全可靠要求！