



# GOTC 2023

## 全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

---

# OPEN SOURCE, INTO THE FUTURE #

---

### 「AI is Everywhere」专场

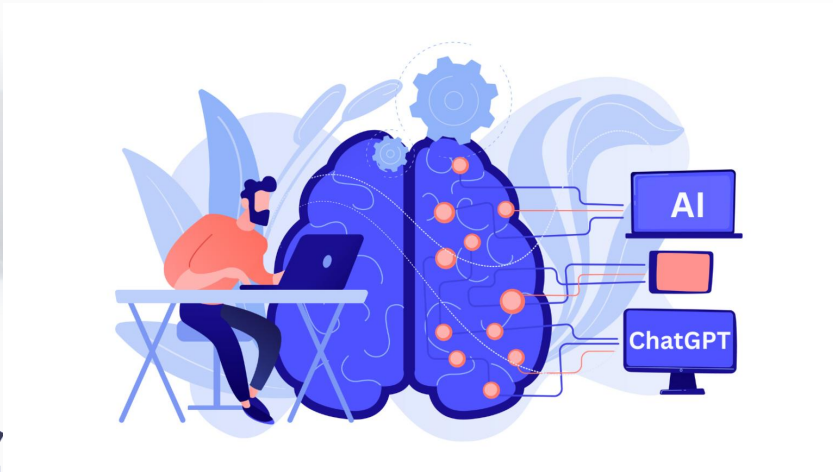
本期议题：旷视算法量产与 MegEngine 生态建设

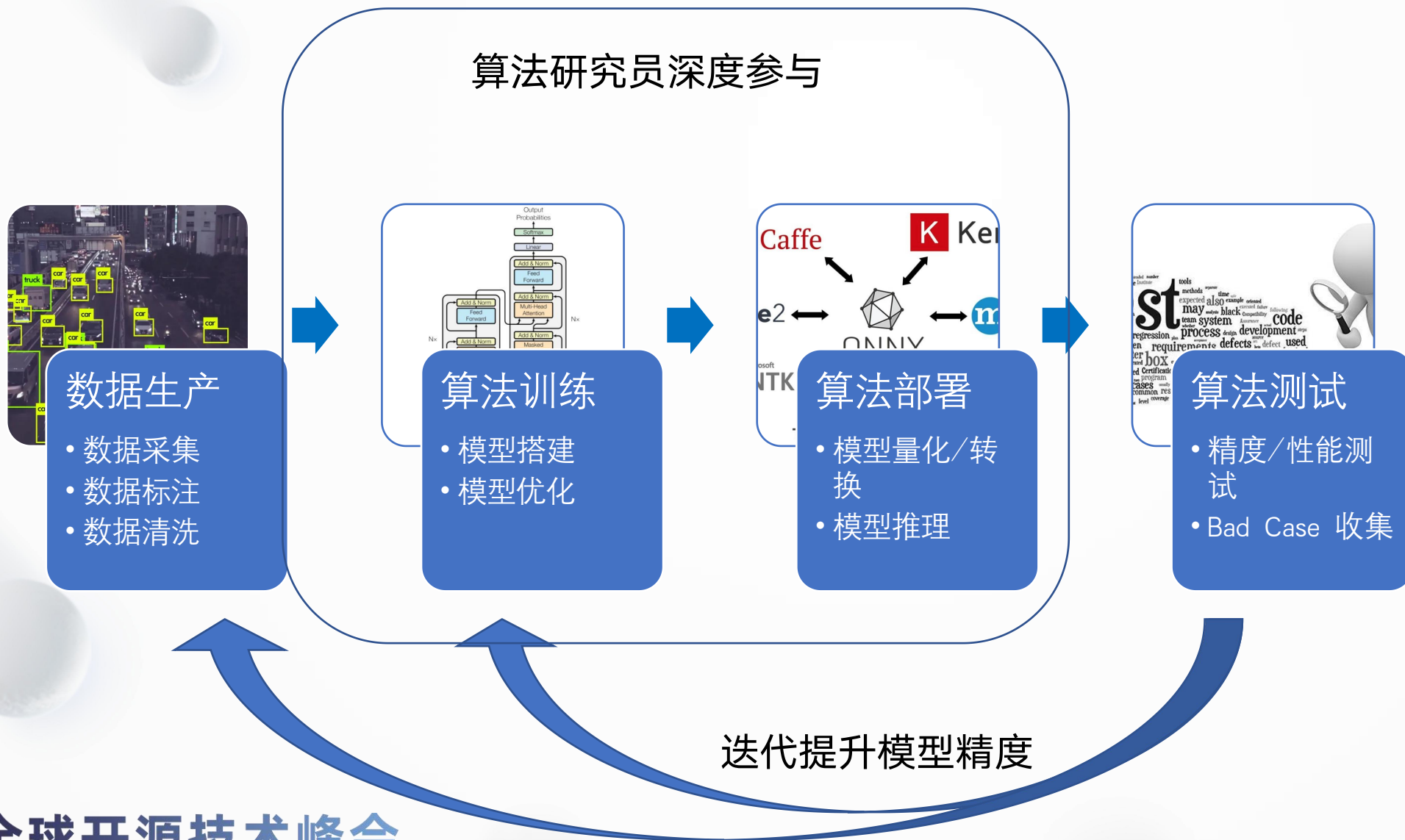
陈其友 2023年6月1日



- 越来越多的领域，如视觉、语音、机器翻译、游戏甚至数据库领域都在使用 AI 算法
- AI 的优势：智能，固定运行时间和内存占用、可移植、容易固化成硬件、响应快速等等

## AI 已经不可替代 AI 需求爆炸式增长



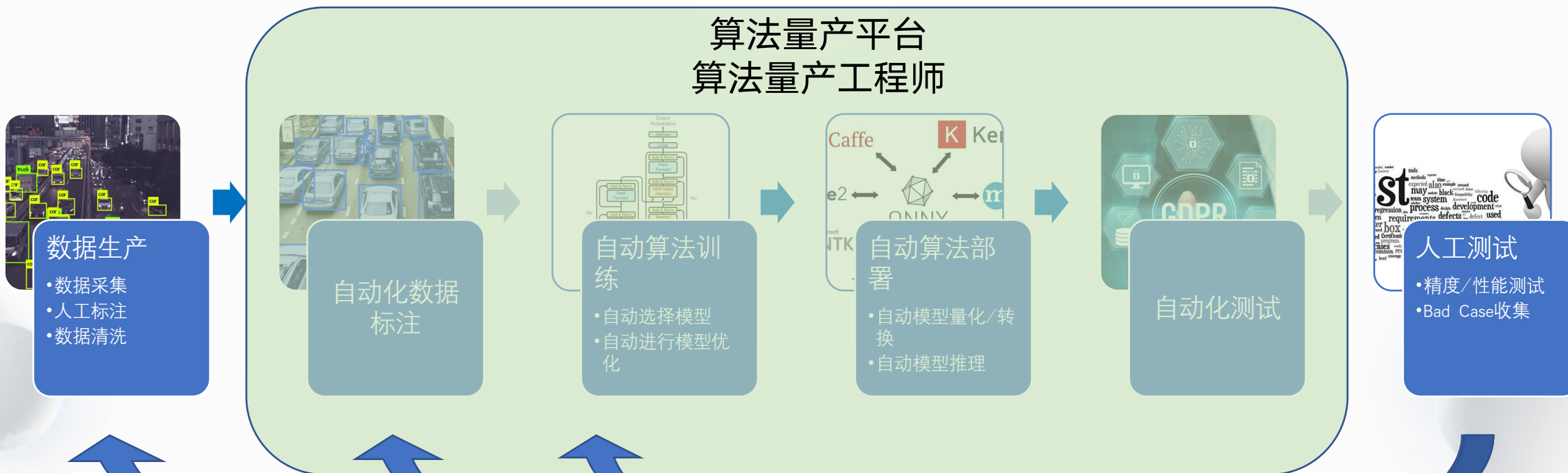
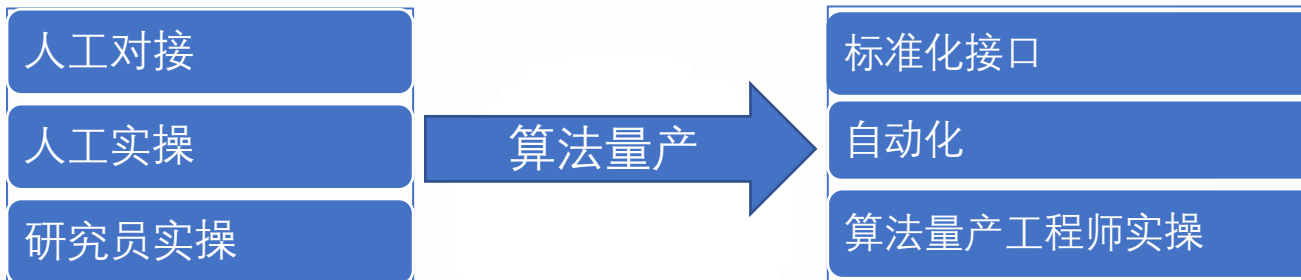


一般算法开发

- 周期长
- 效率低
- 维护成本高
- 高度定制化

算法量产核心思想：

- 流程标准化
- 工作自动化
- 技术低门槛化



迭代提升模型精度

业务端

## 数据管理

数据清洗

数据集管理

数据标注

智能标注

训练对管理

数据质检

## 模型训练

自动学习

实验优化

实验对比

实验管理

workspace

第三方训练

## 模型管理

模型评测

模型转换

## 模型部署

在线部署

服务器算法部署

端设备算法部署

## 运营端

全局用户管理

全局事件管理

资源管理

机器管理

数据统计/审计

## 团队管理

权限管理

事件管理

用户管理

项目管理

## 基础服务

K8S

Nori

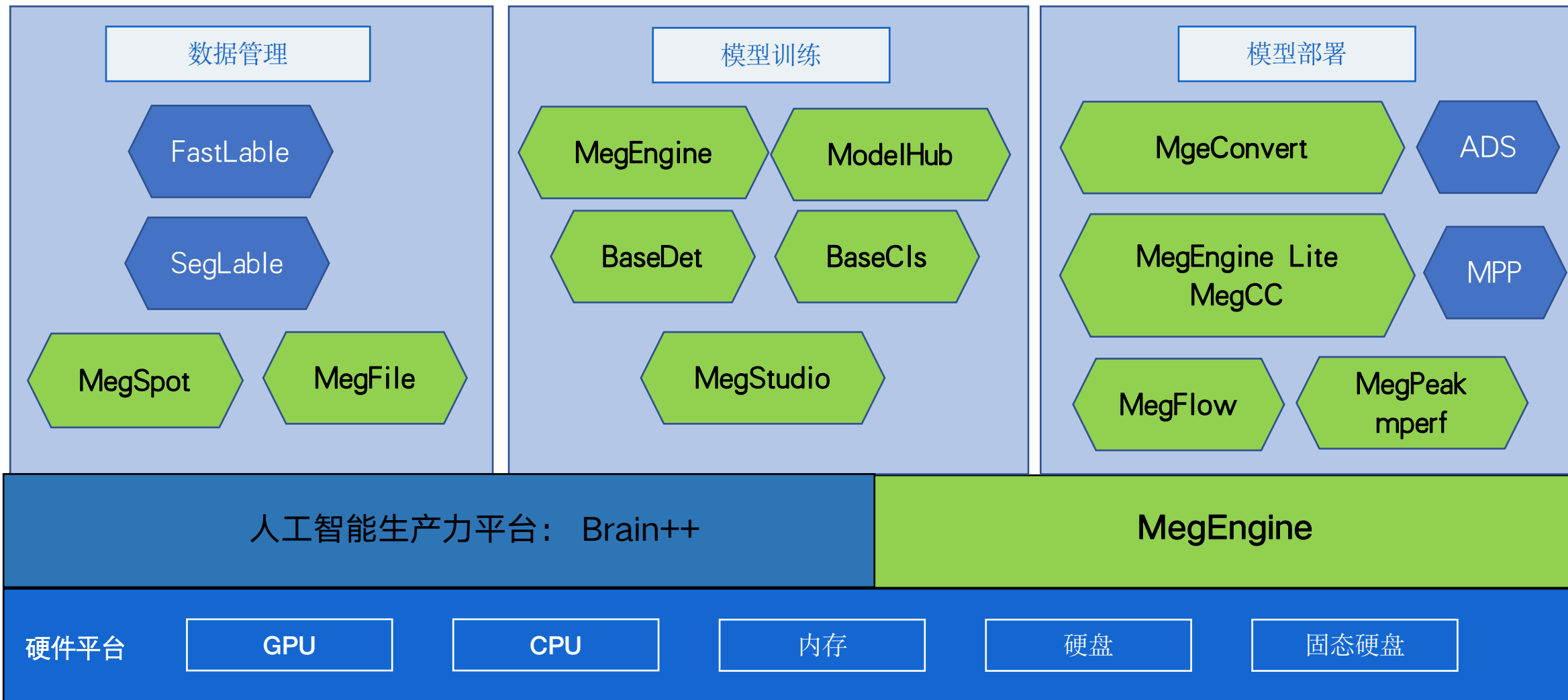
数据库

对象存储

镜像服务

日志服务

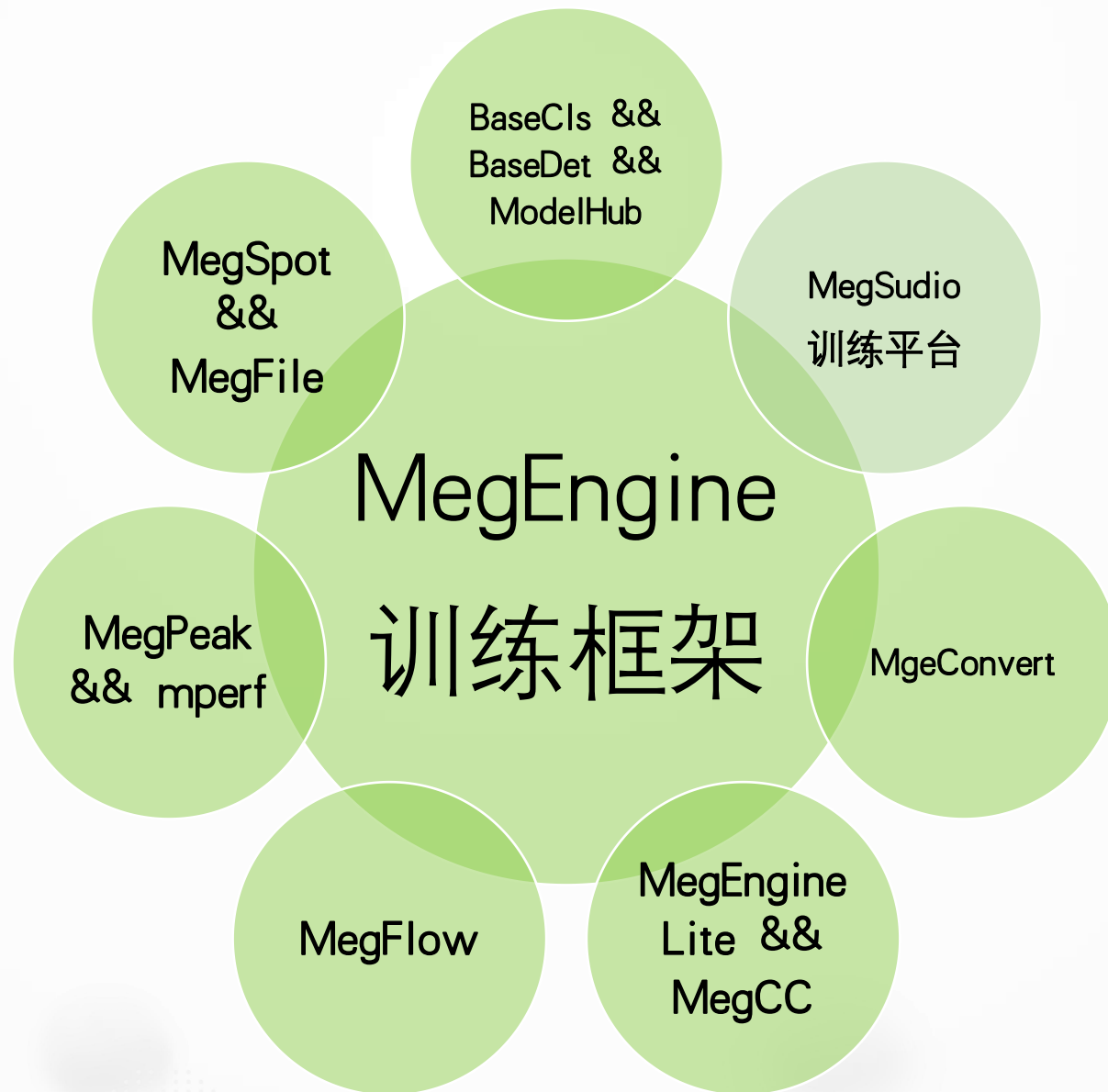
## AIS 旷视 MLOps 算法生产平台





## 覆盖功能

- 数据预处理
- 模型训练
- 模型转换
- 模型推理
- 算法 pipeline 搭建
- 硬件性能测评



## 从 Alpha 到 V1.13，不断优化，提升产品竞争力



- 用户使用体验的优化
- 训练性能，推理性能优化
- 以及生态建设

## 持续迭代，打造自主可控的国产开源框架



## 训练推理一体

一套内核支撑训练到推理，无需模型转换，精度损失最小化。

前后处理可放入计算图推理，训练推理精确对齐，python C++ 不用写两遍。

traced module → megengine lite → megflow，从模型到高并发视频流处理服务只需几行python代码。



## 超低硬件门槛

DTR 算法减少 75% 显存占用，1080 也能训 transformer。

自研 pushdown 内存分配算法，带来最低的内存/显存占用。

自动代码裁剪可使部署文件 binary size 下降 10 倍，有效降低推理硬件成本。

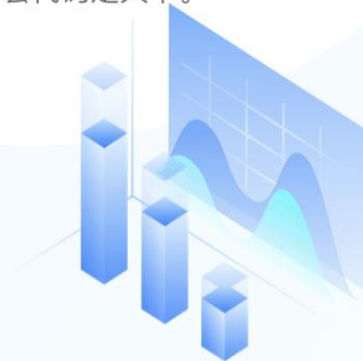


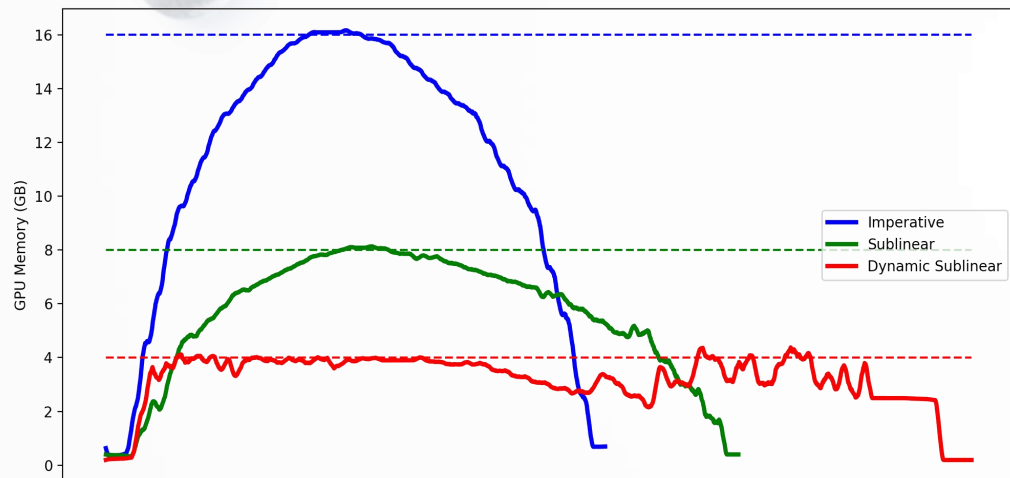
## 全平台高效推理

高效的推理性能，在各类 CPU、GPU 上均可享受到极致性能体验。

自动 layout、kernel 算法选择机制，轻松达到最优推理性能。

跨平台模型精度对齐，一套代码走天下。





```
megengine.dtr.eviction_threshold = "4GB"  
megengine.dtr.enable()
```

训大模型神器

显存占用降至 1/4

动静均支持

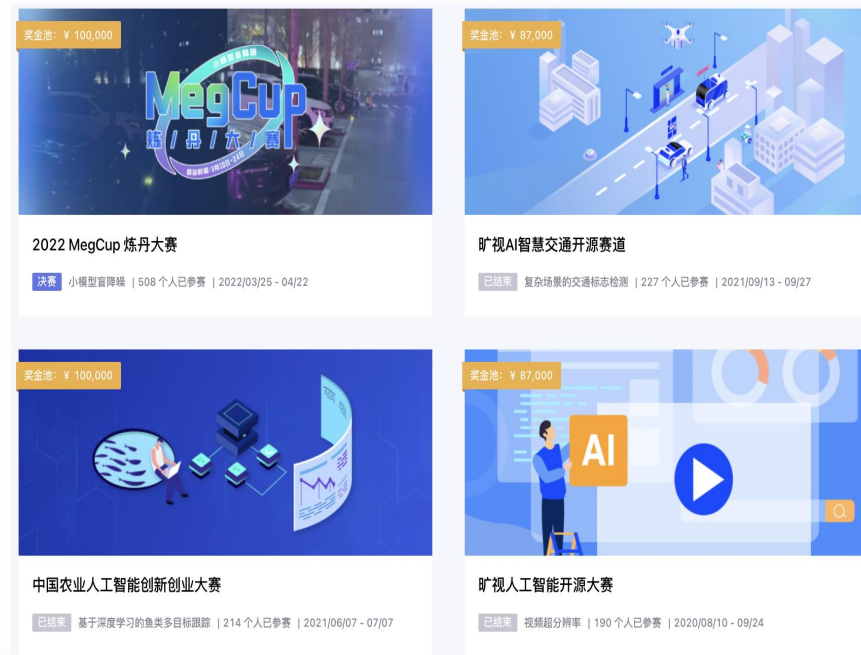
一键开启DTR<sup>[1]</sup>

无需阈值设置

## 可帮助企业 and 高校大幅节省研发硬件投入成本

[1] Kirisame M, Lyubomirsky S, Haan A, et al. Dynamic tensor rematerialization[J]. arXiv preprint arXiv:2006.09616, 2020.

## MegStudio 一站式 “深度学习” 模型开发平台

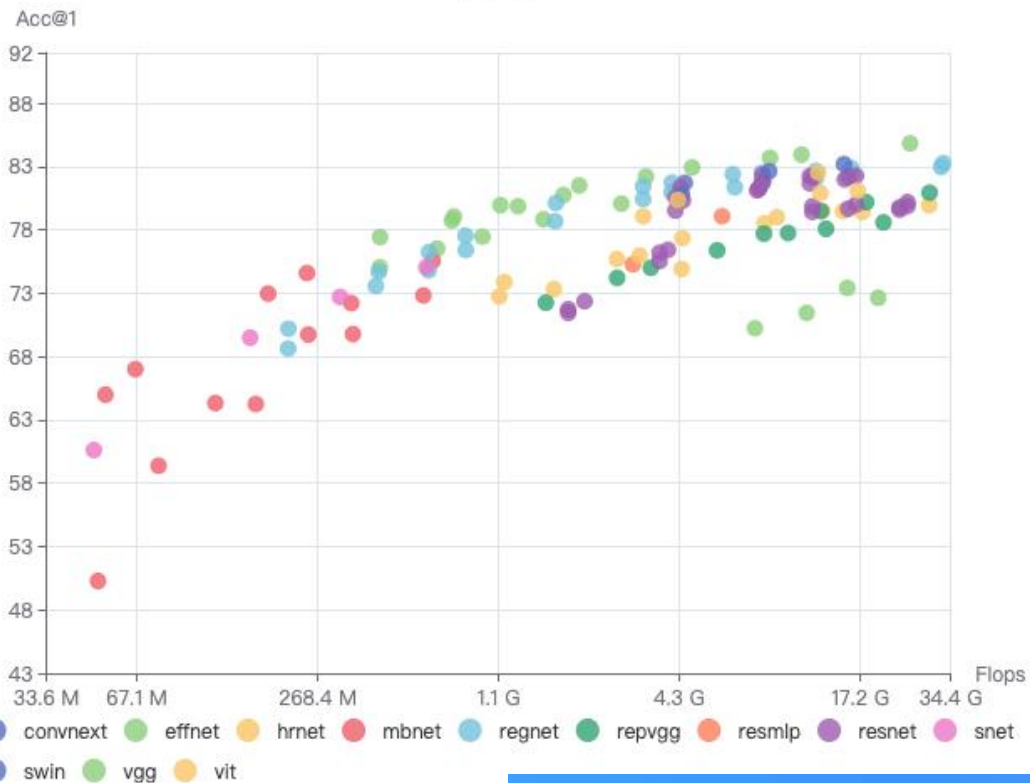


MegStudio 是面向 C 端海量 AI 入门级选手的在线一站式 AI 模型学习与实训平台。为初学者提供了课程、项目、数据集、模型、计算框架、超级算力、持久存储等多种免费资源，大大降低了 AI 技术的准入门槛。



## BaseCls

## BaseDet



### ▼ model zoo

Model name	input size	lr sched	box mAP	weights
Faster R-CNN R50-FPN	800	1x(12e)	37.7	<a href="#">github</a>
RetinaNet R50-FPN	800	1x(12e)	36.2	<a href="#">github</a>
FreeAnchor R50-FPN	800	1x(12e)	38.4	<a href="#">github</a>
FCOS R50-FPN	800	1x(12e)	39.0	<a href="#">github</a>
ATSS R50-FPN	800	1x(12e)	39.5	<a href="#">github</a>
OTA R50-FPN	800	1x(12e)	41.0	<a href="#">github</a>
DETR R50	800	150e	39.9	<a href="#">github</a>

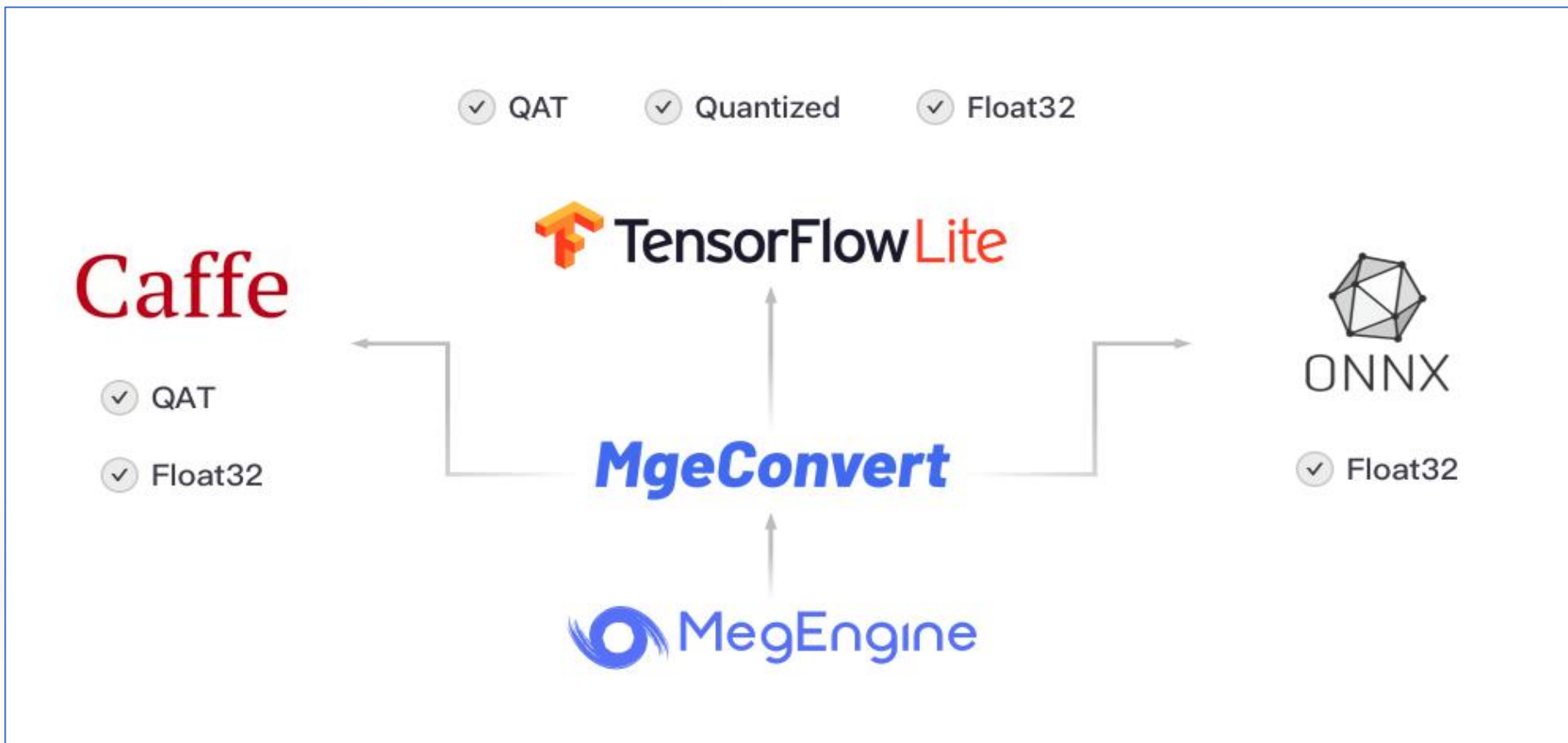
模型中心

[GitHub >>>](#)

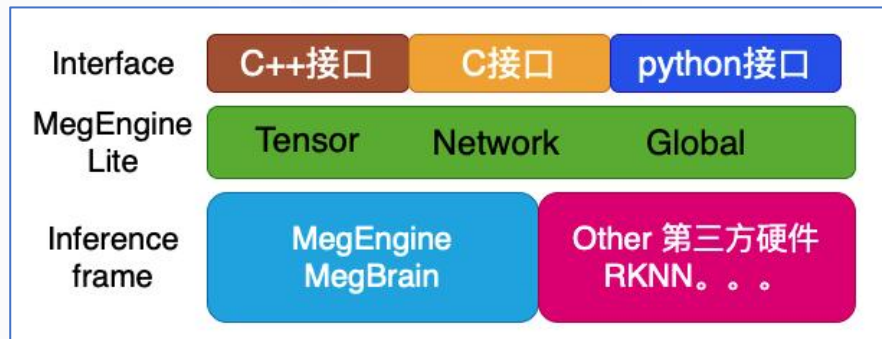
基于旷视研究院领先的深度学习算法，提供满足多业务场景的预训练模型

## MegConvert

MgeConvert 实现多平台模型互转

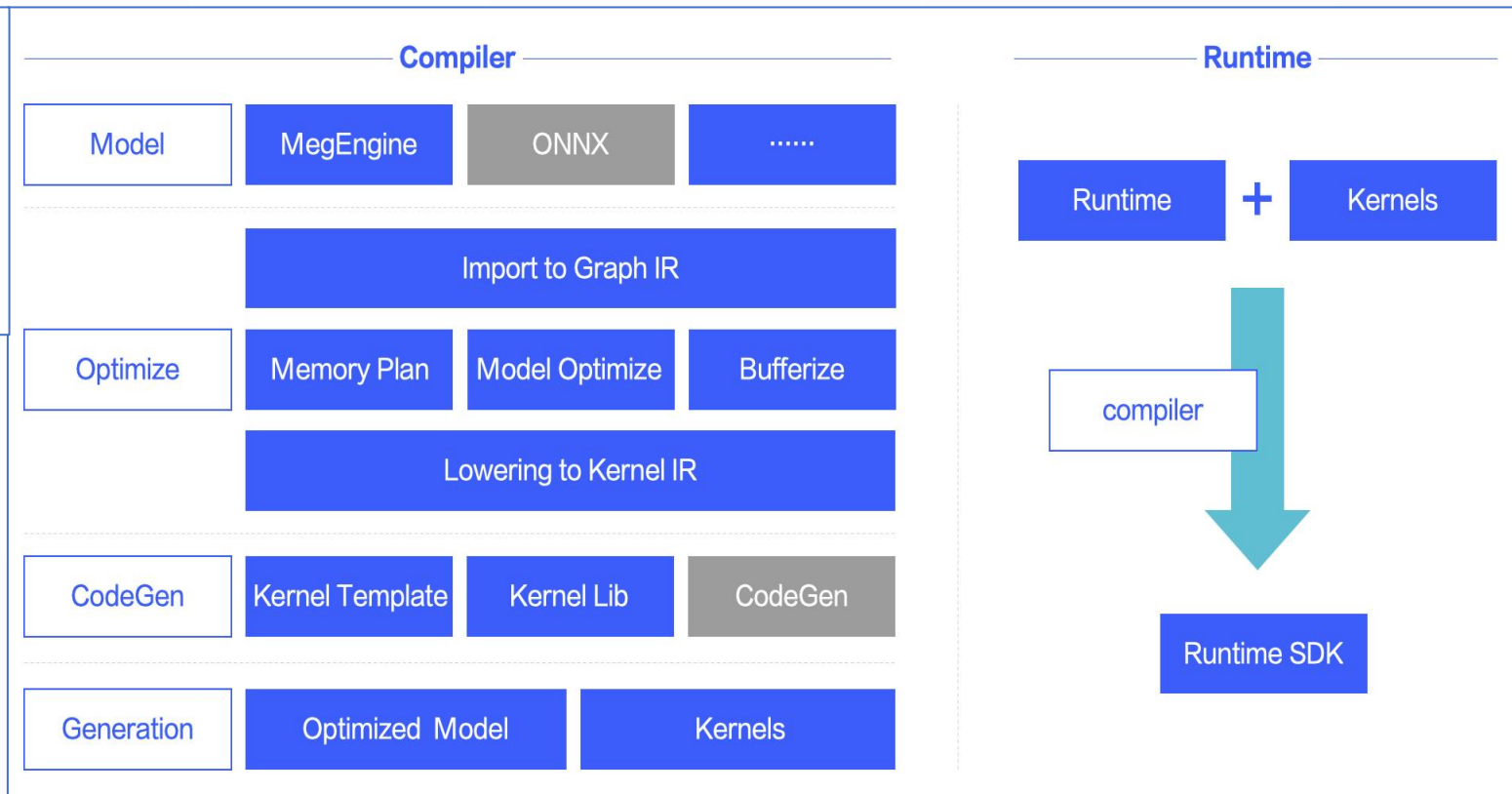


## MegEngine Lite/MegCC



### 覆盖

- 服务器推理
- 移动端推理
- 边缘推理
- 特殊情况推理(MegCC)





## MegFlow



MegFlow 提供了一种更为简洁的视觉应用落地流程，用户可以直接用 Python 搭建计算图，不必关心 C++、图优化相关问题，省去了 SDK 封装流程，可快速实现算法交付。

### 加速 AI 算法交付流程



# MegEngine 生态全景图



## 构建、训练、验证

### MegEngine

使用动态图接口构建、训练和验证你的神经网络，直接使用动态图接口进行推理。

### MegFile

提供完美抽象 S3、HTTP、本地文件等协议的 python 文件接口库。

### MegSpot

提供免费免登录、高效、专业、跨平台的图片&视频的对比的 PC 应用工具。

### ModelHub

基于旷视研究院领先的深度学习算法，提供满足多业务场景的预训练模型。

### BaseCls

极其全面的分类模型库，提供了海量模型的训练代码和预训练权重，全部算法均可以快速部署到硬件上。

### BaseDet

快速且易于使用的检测工具箱，提供了经典的检测 SOTA 模型以及相关组件。

## 平台

### MegStudio

一站式 AI 开发平台，提供框架、算力、编译器、代码托管等。

### MegLab

移动端小程序，提供黑科技 AI 趣味体验，快速感知 AI 魅力。

## 推理、调优、部署

### MegEngine Lite C/C++

利用 Lite C / C++ 接口进行高性能推理。

### MegEngine Lite Rust

利用 rust 接口进行高性能推理。

### MegEngine Lite Python

利用 Lite Python 接口进行更高性能的推理。

### MegCC

MegCC 是一个运行时超轻量、高效、移植简单的深度学习模型编译器。

### MgeConvert

MgeConvert 是 MegEngine 和第三方格式之间实现互联互通的转换工具。

### MegPeak

MegPeak 是一个进行高性能计算的辅助工具，能够使得开发人员轻松获得目标处理器的内在的详细信息，辅助进行对代码的性能评估，以及优化方法设计。

### mperf

mperf 是一个微架构层次的算子性能调优工具箱，主要面向移动/嵌入式平台的 CPU/GPU 核心，目标是“为构建一个更接近闭环的算子调优反馈回路”提供系列基础工具。

### MegFlow

MegFlow 流式计算框架，助力 AI 应用快速落地，简化模型交付流程，实现 15 分钟完成定制化功能。

## 库

- GitHub:  
<https://github.com/MegEngine/MegEngine>
- 官网:  
<https://www.megengine.org.cn/>
- 论坛:  
<https://discuss.megengine.org.cn/>

群号: 1029741705



技术交流 QQ 群

# THANKS