

AI芯片洗牌 演绎哪些新趋势

本报记者 张心怡

无芯片不AI，芯片是支撑人工智能的基础。2019年，云端AI芯片迎来亚马逊、高通、阿里巴巴、Facebook等新玩家，软硬一体化趋势加强；终端芯片功耗比竞争加剧，语音芯片持续火热；边缘AI芯片势头初现。2020年，AI芯片将逐渐进入洗牌期，机遇与挑战并存。



围绕边缘AI芯片的抢滩布局已经开始，头部厂商正在打造云、边、端一体化的计算格局。

边缘AI芯片进入抢滩战

AI正在从云端向边缘端扩展，边缘计算被视为人工智能的下一个战场。寒武纪副总裁刘道福表示，在边缘计算种类中，边缘往往和各类传感器相连，而传感器的数据往往是非结构化的，很难直接用于控制和决策，因此需要边缘人工智能计算将非结构化数据结构化，从而用于控制和决策。

2019年，围绕边缘AI芯片的抢滩布局已经开始。一方面，英伟达、寒

武纪、百度等已经在云、端有所积累的厂商，希望以边缘芯片完善云、边、端生态，打造一体化的计算格局。

英伟达发布了面向嵌入式物联网的边缘计算设备Jetson Nano，适用于入门级网络硬盘录像机、家用机器人以及具备全面分析功能的智能网关等应用，之后英伟达又发布了边缘AI超级计算机Jetson Xavier NX，能够在功耗10W的模式下提供最高14TOPS，在功耗15W模

式下提供最高21TOPS的性能。

寒武纪发布了用于深度学习的SoC边缘加速芯片思元220，采用台积电16nm工艺，最大算力32TOPS (INT4)，功耗控制在10W，支持Tensorflow、Caffe、mxnet以及pytorch等主流编程框架。

百度联合三大运营商、中兴、爱立信、英特尔等企业，发起百度AI边缘计算行动计划，旨在利用AI推理、函数计算、大数据处理和产业模

型训练推动AI场景在边缘计算的算力部署和平台支持。

另一方面，自动驾驶等专用边缘AI芯片势头渐显。地平线宣布量产国内首款车规级AI芯片“征程二代”，采用台积电28nm工艺，可提供超过4TOPS的等效算力，典型功耗控制在2瓦，延迟少于100毫秒，多任务模式下可同时运行超过60个分类任务，每秒识别目标数超过2000个，以应对车联网对强实时响应的需求。

软件是智能操作的核心，随着异构计算逐渐导入AI芯片，软硬件协同成为云端AI的重要趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

高通推出了面向数据中心推理计算的云端AI芯片Cloud AI 100，峰值性能超过350TOPS，相比其他商用方案每瓦性能提升10倍。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

多个新玩家入局云端

云端仍然是AI芯片的主要战场。2019年，云端芯片迎来多个新玩家，算力大战持续升级。

高通推出了面向数据中心推理计算的云端AI芯片Cloud AI 100，峰值性能超过350TOPS，相比其他商用方案每瓦性能提升10倍。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

阿里巴巴推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

功耗比仍是终端侧重点

在终端侧，功耗比仍然是角逐焦点。尤其在手机等对于续航能力锱铢必较的终端，主力厂商推出的AI引擎都对低功耗有所强调。

麒麟990 5G的NPU采用双大核+微核的方式，大核负责性能，

微核拥有超低功耗。据介绍，微核在人脸检测的应用场景下，能耗是大核工作的1/24。高通发布的骁龙865集成了传感器中枢，让终端能够以极低功耗感知周围情境。三星提出通过较低功耗的NPU实现终端设备上的AI处理，

实现在设备端直接执行更复杂的任务。

除了手机，终端侧的另一个当红炸子鸡是AI语音芯片。科大讯飞、阿里巴巴、深鉴科技、清微智能等都发布了针对智能家居的AI语音芯片，反映了AI芯片在特定

领域的专业化、定制化趋势。阿里达摩院公布了首款专用于语音合成算法的AI FPGA芯片技术Ouroboros，使用了端上定制硬件加速技术，降低对云端网络的依赖，支持实时语音合成和AI语音识别，有望率先在天猫精灵搭载。

AI芯片将持续火热，企业扎堆进入。但是2020年前后，将出现一批出局者，行业洗牌开始。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，软硬件协同成为

云端AI的重要趋势。英特尔推出了面向异构计算的统一软件平台One API，以隐藏硬件复杂性，根据系统和硬件自动适配功耗最低、性能最佳的加速方式，简化并优化编程过程。赛灵思也推出了软件平台Vitis AI，向用户开放易于访问的软件接口，可根据软件或算法自动适配赛灵思硬件架构。

AI正在渗入手机和语音芯片，反映了AI芯片在特定领域的专业化、定制化趋势。

2019年，云端芯片迎来多个新玩家，算力大战持续升级。

阿里推出号称“全球最高性能AI推理芯片”含光800，采用自研芯片架构和达摩院算法，在Resnet50基准测试中获得单芯片性能第一。

云服务领跑者亚马逊推出了机器学习推理芯片AWS Inferentia，最

高算力为128TOPS，在AI推理实例infl可搭载16个Inferentia芯片，提供最高2000TOPS算力。

腾讯投资的燧原科技发布了面向云端数据中心的AI加速卡云燧

T10，单卡单精度算力达到20TOPS，支持单精度FP32和半精度BF16的混合精度计算，并为大中型数据中心提供了单节点、单机架、集群三种模式，在集群模式下可通过片间互联实现1024节点集群。

芯片是AI的载体，而软件是完成智能操作的核心。随着异构计算逐渐导入AI芯片，