

## 人工智能芯片：AI 巨轮的引擎



东方证券  
ORIENT SECURITIES

## 报告起因

- 近日，Google 公布其第二代 TPU 产品，Apple 也宣称正在研发一款名为“苹果神经引擎(Apple Neural Engine)”的 AI 专用芯片，人工智能芯片浪潮再起。
- 乌镇围棋峰会，谷歌人工智能 AlphaGo 横扫柯洁等人类顶尖棋手，人工智能再度引发强烈关注。

## 核心观点

- **人工智能芯片前景广阔：**随着下游领域智能化需求的拉动，以及软硬件技术不断取得突破，人工智能技术再次迎来黄金发展期，有望引领新一轮技术革命。作为引领人工智能算法发展方向的深度学习技术，其中的核心环节在于采用 AI 芯片大幅提升计算能力，国内外各大科技巨头纷纷着力布局，人工智能芯片有望实现跨越式增长。
- **各类 AI 芯片百花齐放：**人工智能芯片主要包括 GPU、FPGA、ASIC 以及类脑芯片等类别，在人工智能时代，他们发挥各自优势，呈现百花齐放的态势：GPU 并行计算能力突出，在深度学习训练环节具备优势；FPGA 高性能低功耗的特性适合于推理环节取代传统的 CPU；ASIC 芯片专门针对人工智能设计，有望未来成为主流；类脑芯片作为突破性技术路线，未来如实现突破也将推动人工智能产业长远发展。
- **人工智能芯片在云端与终端领域携手共进：**芯片是人工智能技术的核心环节，当前各大科技巨头在云端 AI 芯片领域进步较快，通过“云端化”+“AI 芯片集群化”的模式高效为用户提供最大化便利；对于终端 AI 芯片领域，目前在汽车、机器人、家居等场景，人工智能技术已经开始得到应用，部分科技巨头也开始切入相应市场进行布局，终端 AI 芯片领域未来有望放量。

## 投资建议与投资标的

- 未来我国人工智能芯片领域有望得到迅猛发展，国内已经有部分企业在沿人工智能产业链进行布局，在核心芯片、大数据、生物识别、物联网、安防等领域，国内公司均已顺利切入并取得一定突破进展。
- 结合公司整体业务和人工智能芯片领域的状况，我们建议关注中科曙光、全志科技、景嘉微、通富微电、富瀚微、海康威视、大华股份。

## 风险提示

- 人工智能芯片研发不及预期；
- 下游需求不及预期

## 行业评级

看好 中性 看淡 (维持)

国家/地区

中国/A 股

行业

电子

报告发布日期

2017 年 05 月 31 日

## 行业表现



资料来源：WIND

## 证券分析师

蒯剑

021-63325888\*8514

kuaijian@orientsec.com.cn

执业证书编号：S0860514050005

胡誉镜

021-63325888\*7518

huyujing@orientsec.com.cn

执业证书编号：S0860514080001

王芳

021-63325888\*6068

wangfang1@orientsec.com.cn

执业证书编号：S0860516100001

## 联系人

王若擎

021-63325888-5023

wangruoqing@orientsec.com.cn

马天翼

021-63325888\*6115

matianyi@orientsec.com.cn

## 目 录

一、	人工智能芯片前景广阔 .....	5
1.1.	人工智能市场高速增长 .....	5
1.2.	深度学习引领人工智能算法发展方向 .....	7
二、	GPU：并行计算能力突出 .....	10
2.1.	GPU 已获得广泛应用 .....	10
2.2.	GPU 的优势来自并行计算能力 .....	11
2.3.	Nvidia 垄断 GPU 市场，国内公司逐步突破 .....	12
三、	FPGA：低功耗场景凸显优势 .....	14
3.1.	FPGA 性能领先 .....	15
3.2.	双寡头垄断 FPGA 市场 .....	17
3.3.	国内 FPGA 产业孜孜求索 .....	19
四、	ASIC：有望成为主流趋势 .....	20
五、	类脑芯片：超越“冯·诺依曼”架构的新思路 .....	22
六、	人工智能芯片在云端与终端携手共进 .....	25
6.1.	云端 AI 芯片领域百家争鸣 .....	25
6.2.	终端 AI 芯片领域初露头角 .....	31
	投资建议 .....	38
	风险提示 .....	40

## 图表目录

图 1：人工智能关键要素.....	5
图 2：全球人工智能市场规模（单位：亿美元） .....	6
图 3：中国人工智能市场规模（单位：亿元） .....	6
图 4：全球人工智能主要公司 .....	6
图 5：国际人工智能领域三巨头动作.....	6
图 6：国内人工智能主要企业 .....	7
图 7：深度学习 VS 神经网络.....	8
图 8：深度学习市场规模.....	8
图 9：深度学习主要市场参与者及开源平台 .....	9
图 10：各公司主要开源平台列表 .....	9
图 11：主要深度学习平台性能比较.....	9
图 12：GPU 在深度学习领域应用广泛 .....	10
图 13：使用 NVidia 加速计算 GPU 的企业数量快速增长 .....	11
图 14：CPU 与 GPU 结构差异 .....	11
图 15：GPU 在 3 年时间内性能提高 50 倍.....	12
图 16：GPU 每秒计算量远超 CPU .....	12
图 17：GPU 是 Nvidia 的主要产品（2016 年报） .....	13
图 18：Nvidia 在 GPU 市场有绝对优势 .....	13
图 19：NVidia 公司加速运算 GPU 及相关产品 .....	13
图 20：Nvidia 近年来财务数据（单位：百万美元） .....	13
图 21：中国在 GPU 领域取得最新成就 .....	14
图 22：FPGA 内部结构原理图 .....	14
图 23：CPU、GPU 及 FPGA 单次迭代时间比较（单位：微秒） .....	15
图 24：CPU、GPU 及 FPGA 单次迭代能耗比较（单位：毫焦） .....	16
图 25：CPU、GPU 及 FPGA 三种芯片性能比较 .....	16
图 26：全球 FPGA 市场规模保持较快增长（单位：亿美元） .....	17
图 27：2016 年 FPGA 市场份额分布.....	17
图 28：英特尔 Lake Crest 架构.....	18
图 29：Canyon Bridge Capital Partners 拟收购 Lattice .....	19
图 30：谷歌 TPU 内部架构.....	21
图 31：寒武纪芯片.....	21
图 32：寒武纪 2 号 DaDianNao 版图 .....	21
图 33：中星微 NPU 架构图 .....	22
图 34：2022 年类脑芯片不同类型终端应用占比.....	23

图 35：各国类脑计算研究项目列表.....	23
图 36：各科技巨头类脑芯片产品列表.....	24
图 37：IBM 第一代 TrueNorth 芯片.....	24
图 38：第一代 IBM TrueNorth 芯片与第二代比较.....	25
图 39：IBM 神经元计算机包含 16 颗 TrueNorth 芯片.....	25
图 40：全球云计算市场规模（亿美元）.....	26
图 41：云计算平台人工智能功能.....	26
图 42：2011 年 Watson 参加节目《Jeopardy》并取得冠军.....	26
图 43：Watson 产生答案流程.....	26
图 44：IBM POWER 处理器发展路径.....	27
图 45：POWER8 架构图.....	27
图 46：微软 Azure 功能.....	28
图 47：2014 年亚马逊 AWS 市场份额占比遥遥领先.....	28
图 48：亚马逊 AWS 能够提供的服务.....	28
图 49：谷歌云计算平台.....	29
图 50：阿里云适用场景.....	30
图 51：阿里云新一代 HPC.....	30
图 52：百度与 Altera 合作建立 FPGA 集群.....	31
图 53：百度开放云功能.....	31
图 54：Nvidia Drive PX 车载计算平台.....	32
图 55：Nvidia Drive PX2 平台.....	32
图 56：Nvidia Xavier 芯片.....	33
图 57：高通发布智能汽车芯片 602A.....	33
图 58：国内汽车电子芯片市场规模.....	34
图 59：飞思卡尔 Vybrid 处理器.....	35
图 60：赛灵思 FPGA 芯片.....	35
图 61：夏普机器人手机 RoBoHoN.....	35
图 62：亚马逊 Echo 音箱基本构造.....	36
图 63：Echo 音箱主板芯片构成.....	36
图 64：京东&科大讯飞叮咚音箱.....	37
图 65：叮咚音箱主板构造.....	37
图 66：人工智能芯片及应用.....	37
图 67：A 股上市公司切入人工智能领域情况.....	38

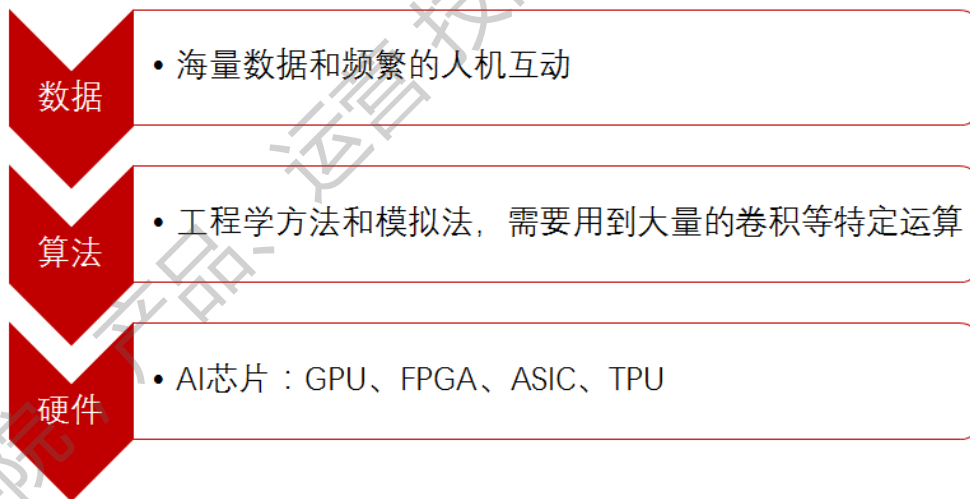
## 一、 人工智能芯片前景广阔

人工智能（AI, Artificial Intelligence）是用于开发和研究用于模拟甚至扩展人的智能的技术及应用系统的一门新的技术科学。人工智能的目标是对人意识和思维过程的模拟，让机器做到像人一样思考，甚至超过人的智能，从而使机器能够胜任通常需要人类智能才能完成的复杂工作。

当前实现人工智能的主要途径是软件算法。目前算法主要可以分为工程学方法和模拟法两种，工程学方法利用大量数据处理经验，运用传统的编程技术使系统呈现智能效果，该方法已经在文字识别等领域有所建树；模拟法则在运算结果和实现方法两个维度模仿人类或其他生物机理，从而提升算法性能，遗传算法（GA）及神经网络（ANN）均属于此类算法。人工智能算法不同于常规算法，需要用到大量的卷积等特定运算，常规处理器芯片在进行这些运算时效率较低，人工智能算法需要特殊的芯片。

目前主流芯片为 GPU 并行计算神经网络，而 FPGA 和 ASIC 也将成为推动人工智能进步的强大动力。

图 1：人工智能关键要素



数据来源：艾瑞咨询，东方证券研究所

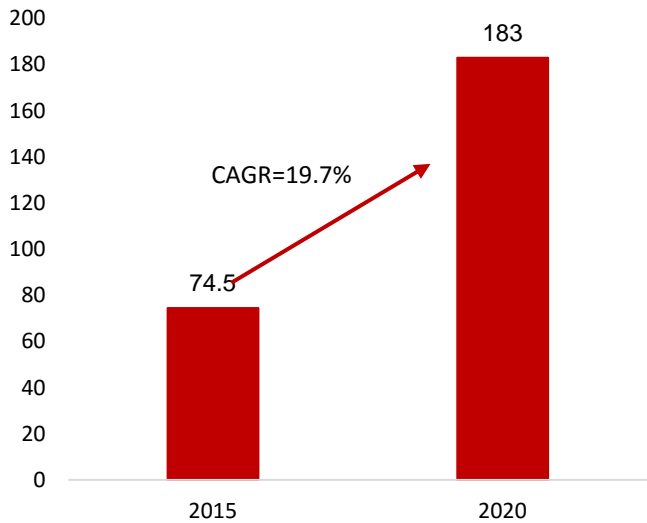
### 1.1. 人工智能市场高速增长

在人工智能超过 60 年的发展历程中，经历了漫长的历史演进和技术更迭，并曾两次陷入低谷。近几年随着工业 4.0、智能生活、“互联网+”等领域的快速进步，加之深度学习算法在语音和视觉识别上取得突破，人工智能技术开始渗透至工业、医疗、教育、安全等多个领域，尤其是近两年来，由 DeepMind 公司开发的人工智能机器人 AlphaGo 接连击败李世石、柯洁等著名围棋选手，人工智能受到了全球大范围关注，迎来了第三个黄金发展时期。

根据艾瑞咨询的报告，2015 年全球人工智能市场规模为 74.5 亿美元，而到 2020 年市场规模将扩大至 183 亿美元，复合年增长率将达到 19.7%。同时预计到 2020 年，中国人工智能市场规模将从 2015 年的 12 亿元增长至 91 亿元人民币，复合年增长率将达到 50.0%

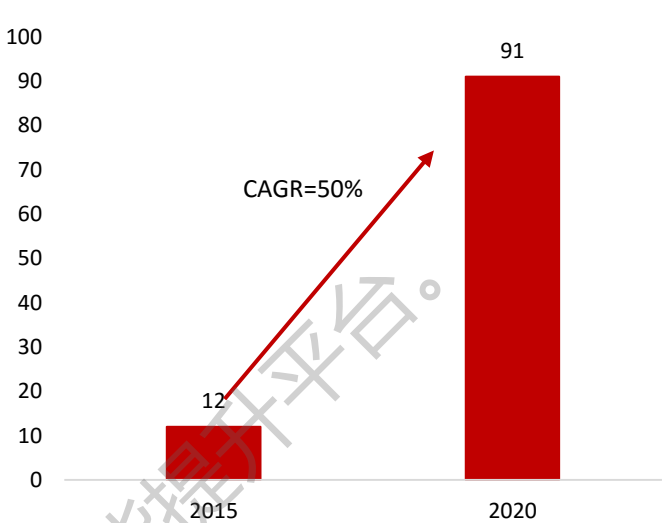


图 2：全球人工智能市场规模（单位：亿美元）



数据来源：艾瑞咨询，东方证券研究所

图 3：中国人工智能市场规模（单位：亿元）

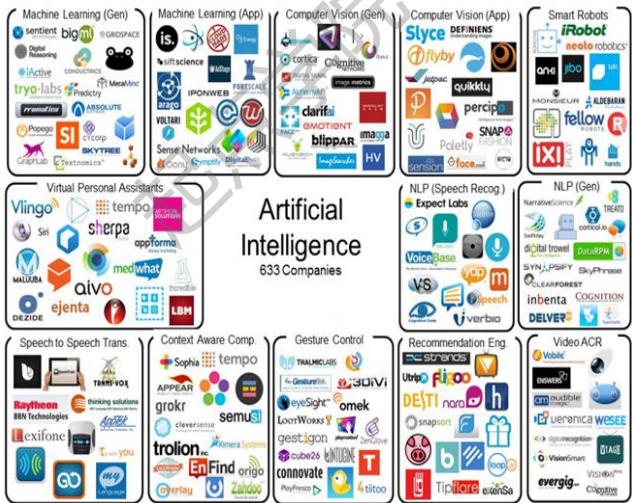


数据来源：艾瑞咨询，东方证券研究所

人工智能市场规模的快速增长得益于其应用领域的不断开拓。根据 Venture Scanner 对全球 957 家人工智能公司的跟踪调查，目前已经覆盖包括深度学习、机器视觉、指纹识别、人脸识别、个人助理、智慧机器人等 13 个具体应用，涉及工业机器人、安全识别、无人驾驶、智能医疗、智能家居等多个新兴产业，人工智能势必将成为新一轮科技革命的强大推动力量。

正因为此，国际科技公司巨头正在加速在人工智能领域的布局。谷歌、微软和英特尔等公司均在该领域不断深耕，取得巨大进展。

图 4：全球人工智能主要公司



数据来源：Venture Scanner，东方证券研究所

图 5：国际人工智能领域三巨头动作

公司	进展	应用领域
谷歌	推出基于人工智能的新搜索算法 RankBrain	智能搜索
	联合福特研发无人驾驶汽车	无人驾驶
	开源人工智能平台 TensorFlow	深度学习
	推出基于人工智能的聊天软件	智能机器人
微软	推出第三代“微软小冰”	智能机器人
	开源机器学习工具包 DMTK	机器学习
	推出人脸情绪识别器	人脸识别
	人工智能助理小娜登陆各个平台	智能机器人
英特尔	6000 万美元投资无人机公司 Yuneec	无人机
	5000 万美元投资量子计算机	硬件升级
	收购人工智能公司 Saffron	数据挖掘
	167 亿美元收购 Altera	芯片制造

数据来源：互联网，东方证券研究所

在国内市场，百度、科大讯飞、阿里巴巴、腾讯等巨头也纷纷在人工智能领域着力布局，而人工智能的广阔前景也吸引国内上百家创业公司投入其中，主要聚焦领域包括智能语音、机器视觉、数据挖掘、智能机器人、无人机等。

图 6：国内人工智能主要企业



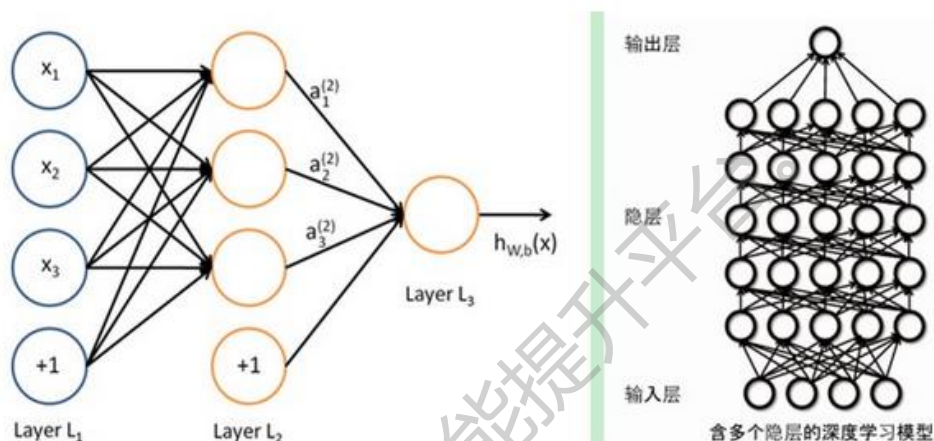
数据来源：艾瑞咨询，东方证券研究所

## 1.2. 深度学习引领人工智能算法发展方向

目前深度学习作为人工智能最主流的算法获得广泛关注。这一概念由 Hinton 等人于 2006 年提出，其实质是通过构建具有很多隐层的机器学习模型和海量的训练数据，使机器去学习更有用的特征，从而最终提升分类或预测的准确性。也就是说，深度学习是对不同模式进行建模的一种方式，其结构具有较多层数的隐层节点以保证模型的深度；同时深度学习明确突出了特征学习的重要性，其通过逐层特征变换，将样本在原空间的特征表示变换到一个新特征空间，从而使识别或预测更加准确。

因此深度学习集中体现了机器学习算法的三大趋势，首先是用较为复杂的模型降低模型偏差，二是用大数据提升统计估计的准确性，三是用可扩展的梯度下降算法求解大规模优化问题。

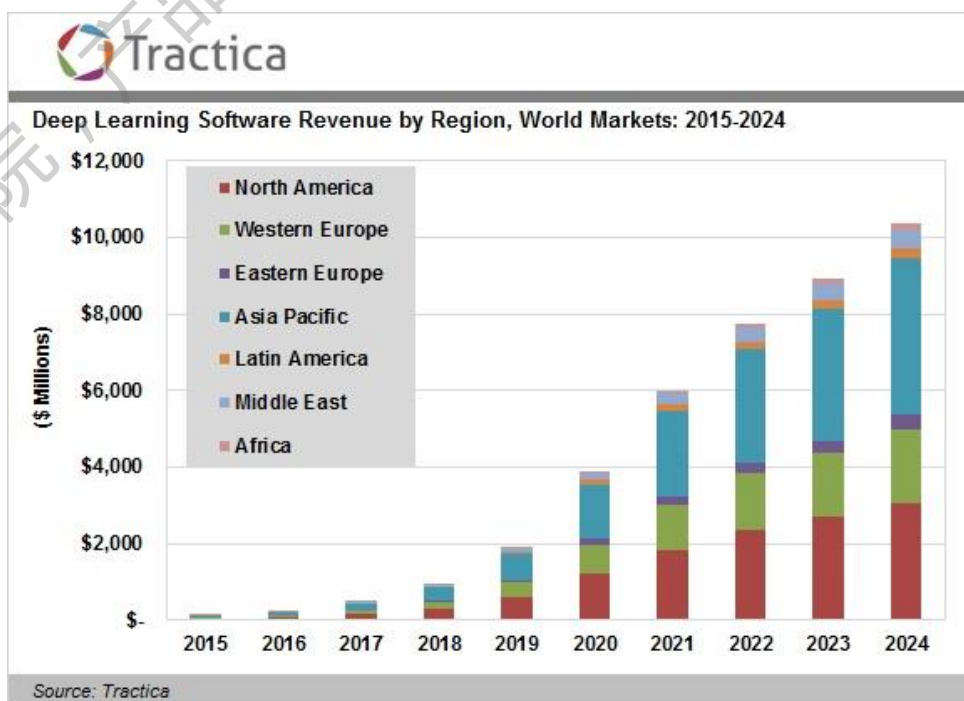
图 7：深度学习 VS 神经网络



数据来源：互联网，东方证券研究所

目前“大数据+深度神经网络”模型已经成为机器学习发展的核心路径，根据 Tractica 的预测，到 2024 年，深度学习仅仅在软件方面的市场价值就将超过 104 亿美元，硬件和服务方面的收入将会是软件市场规模的数倍以上。

图 8：深度学习市场规模



数据来源：Tractica，东方证券研究所

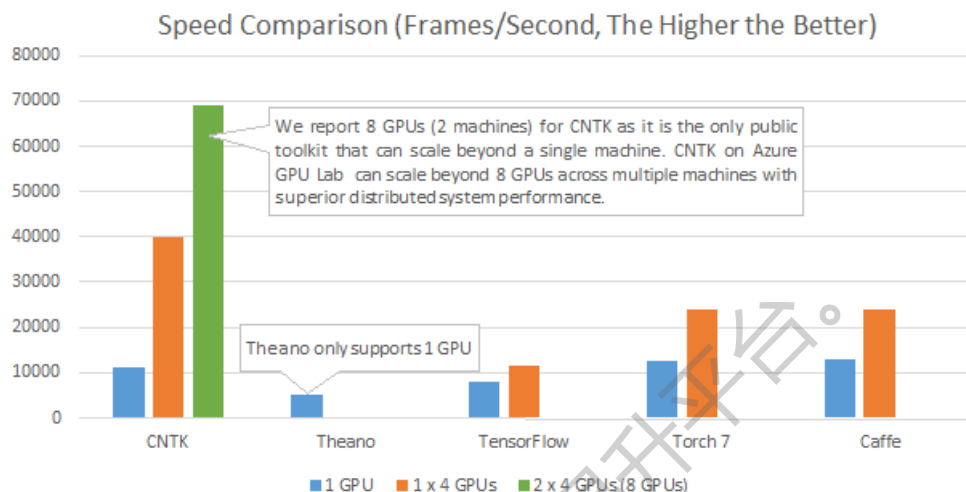


图 9：深度学习主要市场参与者及开源平台



数据来源：各公司官网、东方证券研究所

图 11: 主要深度学习平台性能比较



数据来源：微软公司官网，东方证券研究所

深度学习的兴起得益于大数据的发展、计算机计算能力的大幅提升和算法本身的突破，其中计算能力的大幅提升则得益于 GPU、FPGA、ASIC 等人工智能芯片的广泛应用，芯片作为人工智能技术核心环节，未来前景广阔。

## 二、 GPU：并行计算能力突出

### 2.1. GPU 已获得广泛应用

GPU 即图形处理器，原本是在个人电脑、工作站、游戏机和一些移动设备上专门进行图像运算工作的微处理器。由于其强大的并行计算能力，GPU 逐渐成为目前深度学习领域使用最为广泛的核心芯片。

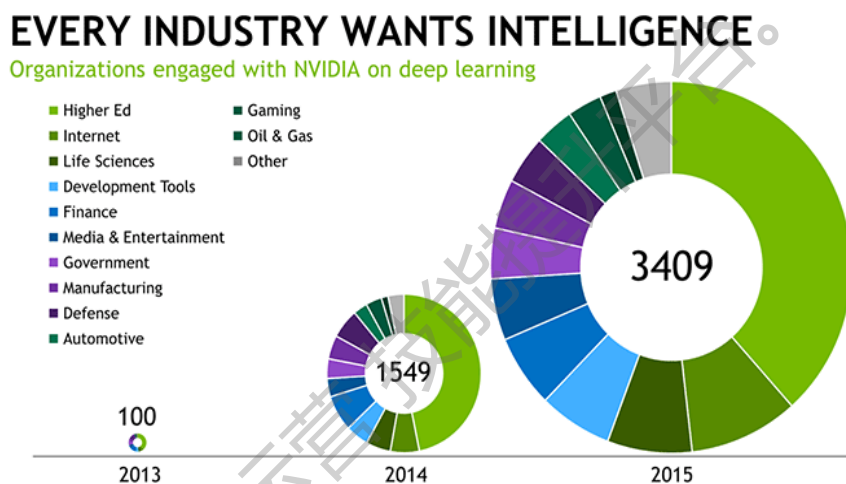
图 12：GPU 在深度学习领域应用广泛



数据来源：TechTarget，东方证券研究所

GPU 已经在图像识别、人脸识别、语音识别、视频分析、自然语言处理等多个领域大放异彩，并逐渐向医药、安全、能源等领域渗透。下游应用的不断扩展反过来又催生了加速计算 GPU 的快速发展。

图 13：使用 NVidia 加速计算 GPU 的企业数量快速增长



数据来源：Nvidia，东方证券研究所

## 2.2. GPU 的优势来自并行计算能力

GPU 与 CPU 有相同之处，两者都有总线和外界联系，都有自己的缓存系统，以及数字和逻辑运算单元。

两者也具有很大的差异。CPU 需要很强的通用性来处理各种不同的数据类型，同时又需要进行逻辑判断、分支跳转和中断等处理，因此 CPU 内部的结构异常复杂；而 GPU 专门执行复杂的数学几何计算，面对的是类型高度统一、相互无依赖的大规模数据和不需要被打断的纯净计算环境。

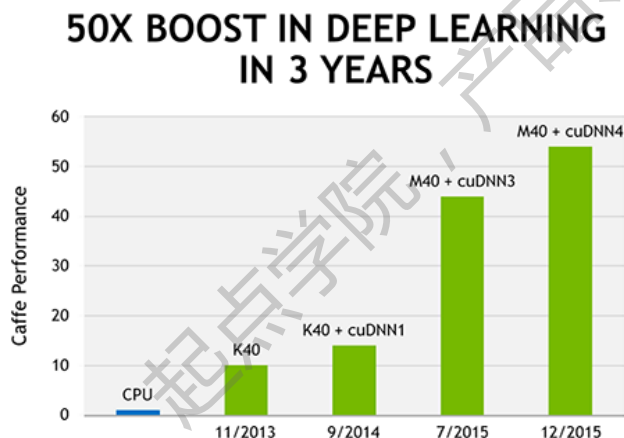
图 14：CPU 与 GPU 结构差异



数据来源：Nvidia CUDA，东方证券研究所

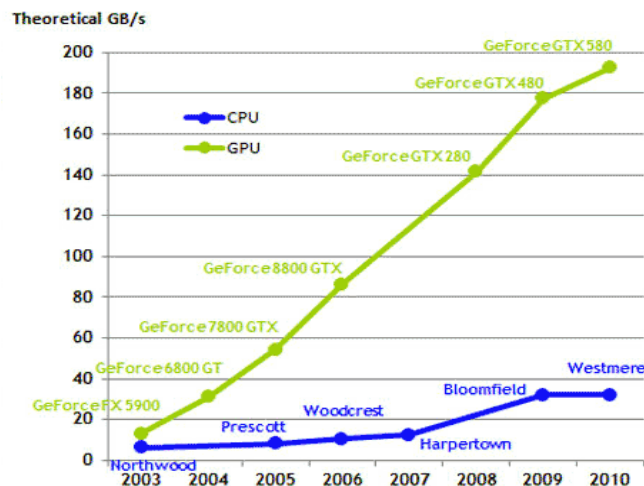
设计目的和处理数据方式的不同导致 GPU 和 CPU 在设计结构上有着天壤之别，包括：片内缓存体系和数字逻辑运算单元结构。CPU 不断增加处理器中晶体管的数量，但在运行单线程串行程序过程中，这些晶体管大多数被用作组成高速缓存，这样做虽然把处理器的功耗控制在合理范围内，但也阻碍了性能的进一步提高；GPU 采用数量众多的计算单元和超长的流水线，但只有非常简单的控制逻辑而省去了高速缓存，所以与 CPU 擅长逻辑控制和通用类型数据运算不同，GPU 擅长大规模、独立的浮点和并行计算，例如计算机图像处理。

图 15：GPU 在 3 年时间内性能提高 50 倍



数据来源：东方证券研究所

图 16：GPU 每秒计算量远超 CPU



数据来源：AllegroViva，东方证券研究所

基于深度学习需要在成千上万的变量中寻找最佳值，并不断通过尝试实现收敛的特性，GPU 自身具备的高并行度、矩阵预算和强大的浮点计算能力可以大幅加速深度学习模型的训练，在相同精度下能提供更快的处理速度、更少的服务器投入和更低的功耗，成为开启深度学习大门的重要推手。

### 2.3. Nvidia 垄断 GPU 市场，国内公司逐步突破

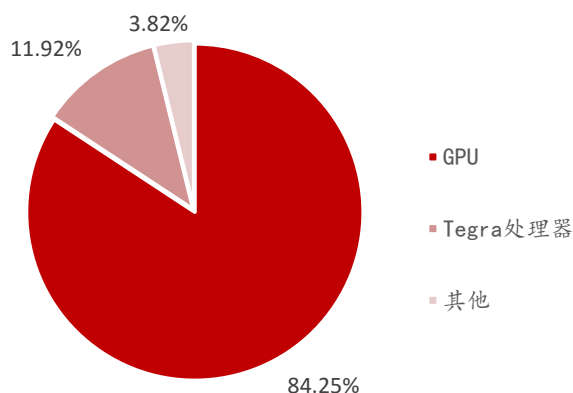
有关分析师的申明，见本报告最后部分。其他重要信息披露见分析师申明之后部分，或请与您的投资代表联系。并请阅读本证券研究报告最后一页的免责申明。

起点学院 www.qidianla.com，人人都是产品经理旗下品牌，打造最专业最系统的产品、运营课程

点击进入 <http://www.hibor.com.cn>

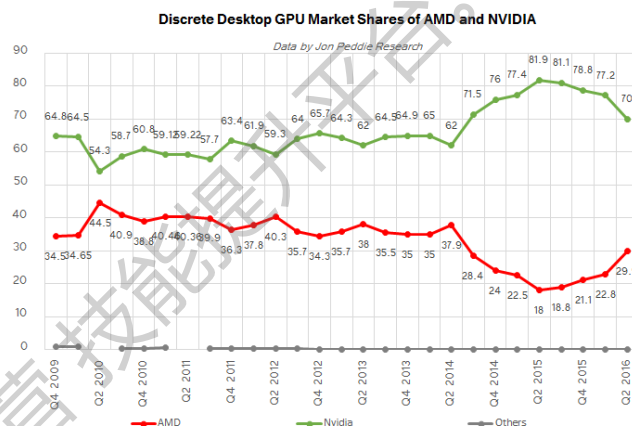
Nvidia 作为 GPU 领域当之无愧的王者，在视觉计算方面拥有数以千计的专利发明和突破性技术。GPU 作为其核心产品，占据其 84% 的收入份额，应用领域涵盖视频游戏、电影制作、产品设计、医学诊断以及科学研究等各个门类。Nvidia 很早就开始深度神经网络的研究并致力于开发加速运算 GPU，目前 Nvidia 已经与谷歌、微软、IBM、丰田、百度等诸多尝试利用深度神经网络来解决海量复杂计算问题的企业建立合作关系，近年来，公司 GPU 出货量的市场份额维持在 70% 以上的绝对优势地位，远远超过 AMD 等竞争对手。

图 17: GPU 是 Nvidia 的主要产品 (2016 年报)



数据来源：Wind，东方证券研究所

图 18: Nvidia 在 GPU 市场有绝对优势



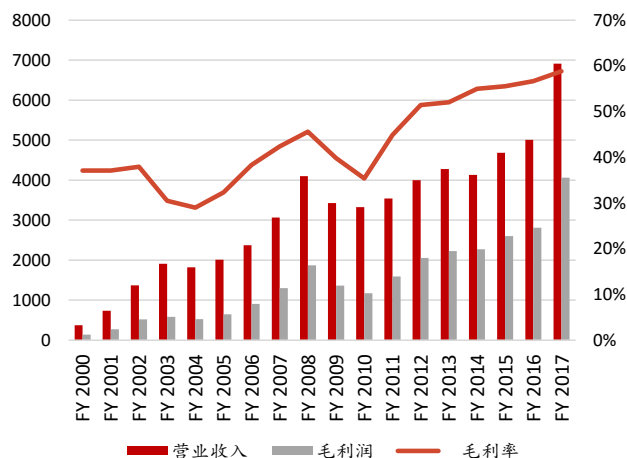
数据来源：WCCF Tech，东方证券研究所

Nvidia 与下游客户在深度学习领域的合作不断加深，已经开发出多款为深度学习量身打造的 GPU 产品，优势的市场地位使其过去几年的毛利率维持在 50% 以上的较高水平。在今年的 GTC 大会上，NVIDIA CEO 黄仁勋发布了首款 Volta 架构的 GPU——GV100 以及产品——Tesla V100 加速卡，Volta 是一款全新的架构，采用了台积电 12nm FFN 制程，相较于之前的 Pascal 架构的 GPU 产品是一次质的飞跃。

图 19: NVidia 公司加速运算 GPU 及相关产品

时间	产品	性能
2015年4月	GeForce GTX TITAN X	全球最快的GPU，采用Nvidia Maxwell GPU架构的TITAN X，结合3072个处理核心，单精度峰值性能为7teraflops，12GB显存，336.5GB/S带宽
2015年4月	DIGITS DevBox平台	采用四个TITAN X GPU，包含DIGITS软件包，和完整的GPU加速深度学习库cuDNN2.0
2015年4月	DRIVE PX	用于自动驾驶汽车的深度学习平台，定位是自动驾驶车载电脑。基于Tegra X1s处理器
2015年4月	Pascal架构	混合精度计算使GPU能在16位浮点精度下拥有两倍于32位浮点精度下和容量的计算速度；采用3D堆叠显存提高近三倍带宽；传输速度将是目前PCI-Express标准的5-12倍
2016年4月	GP100	包含153亿个晶体管芯片，CUDA核心增加至1792个，使精度提升至5.3TFLOPS
2017年5月	GV100	211亿晶体管，815平方毫米，基于台积电的12nm FFN制程

图 20: Nvidia 近年来财务数据 (单位: 百万美元)





数据来源：Bloomberg，东方证券研究所

数据来源：Bloomberg，东方证券研究所

中国在 GPU 芯片设计领域发展相对较晚，当前掌握核心技术的公司包括景嘉微、兆芯等。其中景嘉微研发的 JM5400 图形芯片打破国外芯片在我国军用 GPU 领域的垄断，实现了军用 GPU 国产化。

图 21：中国在 GPU 领域取得最新成就

时间	公司	产品	进展
2014 年	景嘉微	JM5400	与龙芯为合作伙伴，芯片主要应用在军用飞机和神舟飞船上
2016 年	兆芯	ZX-2000	公司主要技术来源于台湾威盛授权，图形核心为美国 S3 Graphics

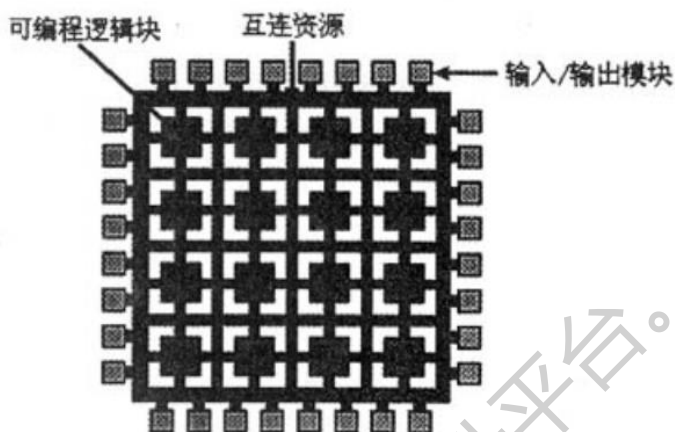
数据来源：互联网，东方证券研究所

### 三、FPGA：低功耗场景凸显优势

FPGA，即现场可编程门阵列，它是在 PAL、GAL、CPLD 等可编程器件的基础上进一步发展的产物，并作为专用集成电路(ASIC)领域中的一种半定制电路而出现，主要为了解决 ASIC 由于大规模工业化生产而导致的结构固化，无法满足某些特定逻辑结构要求的弊端。

FPGA 主要由三部分构成：可配置逻辑模块 (CLB)、输出输入模块 (IOB) 和内部连线 (Interconnect)。可编程逻辑块 (CLB) 是 FPGA 的主要组成部分，是实现逻辑功能的基本单元，可以根据设计灵活地改变连接和设置，完成不同的逻辑功能；输入/输出模块 (IOB) 是芯片和外界接口，提供器件引脚和内部逻辑阵列之间的连接，完成不同电器特性下的输入/输出功能；内部连线 (Interconnect) 包括各种长度的金属连线线段和一些可编程连接开关，它们将各个 CLB 之间以及 CLB 与 IOB 之间互相连接起来，构成各种复杂功能的系统。

图 22：FPGA 内部结构原理图



数据来源：《基于 FPGA 的图像处理算法的研究与硬件设计》，东方证券研究所

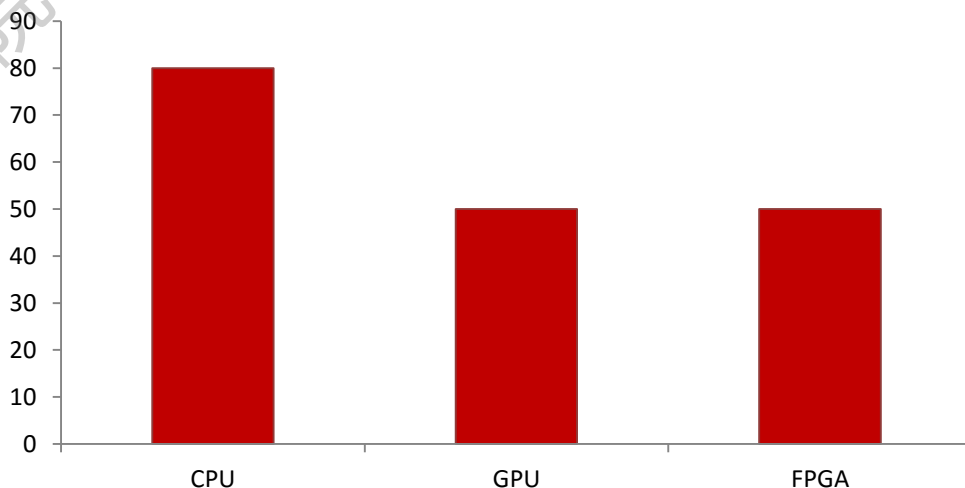
### 3.1. FPGA 性能领先

FPGA 与 GPU 以及 CPU 相比，具有性能高、能耗低以及可硬件编程的特点。

虽然 FPGA 的频率一般比 CPU 低，但是可以用 FPGA 实现并行度很大的硬件计算器。比如一般 CPU 每次只能处理 4 到 8 个指令，在 FPGA 上使用数据并行的方法可以每次处理 256 个或者更多的指令，因此 FPGA 的数据吞吐量远超 CPU。

根据微软研究院对 CPU、GPU 及 FPGA 在加速计算方面的研究，FPGA 和 GPU 算法的单次迭代时间均优于 CPU，且随着矩阵运算规模的增加，GPU 与 FPGA 相比于 CPU 的加速优势会越来越明显。

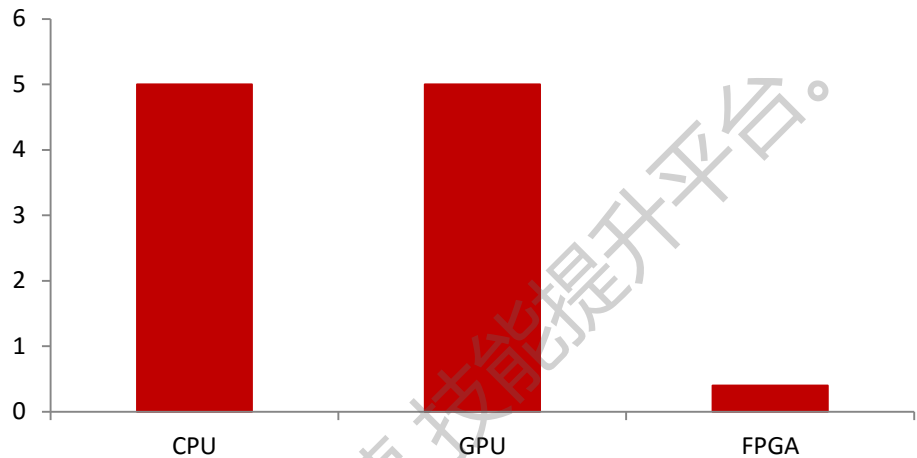
图 23：CPU、GPU 及 FPGA 单次迭代时间比较（单位：微秒）



数据来源：和讯名家，东方证券研究所

并且，FPGA 在能耗方面具有明显的优势。CPU 的解码器通常会占总能耗的 50%，而在 GPU 中，即使其解码器的部分相对较小，也会消耗 10%-20% 的能源。相比之下，由于 FPGA 内部结构没有解码器，加之 FPGA 的主频比 CPU 及 GPU 低很多，通常 CPU 和 GPU 的主频在 1-3GHz 之间，而 FPGA 的主频在 500MHz 以下，因此，FPGA 的能耗要远低于 CPU 及 GPU。

图 24：CPU、GPU 及 FPGA 单次迭代能耗比较（单位：毫焦）



数据来源：和讯名家，东方证券研究所

FPGA 支持硬件编程。FPGA 能够使用户较为方便的设计出所需的硬件逻辑，而且可以进行静态重复编程和动态系统重配置，使系统的硬件功能可以向软件一样通过编程来修改，实现灵活而方便的更新和开发，大大提高系统设计的灵活性和通用性。

图 25：CPU、GPU 及 FPGA 三种芯片性能比较

硬件	CPU	GPU	FPGA
单次迭代时间（微秒）	80	50	50
单次迭代能耗（毫焦）	5	5	0.4
开发难度	小	较小	大
增加功能	容易	容易	难
硬件升级	无需修改代码	无需修改代码	需要修改代码
性能/成本	高	低	高
片外存储器	内存，容量大，速度低	显存，速度高，容量大	内存，速度低
开发周期	短	短	长

注：使用芯片产品为微软芯片基于 BLAS 算法 FPGA 和 GaxPy 算法 CPU、GPU。

数据来源：微软官网，东方证券研究所

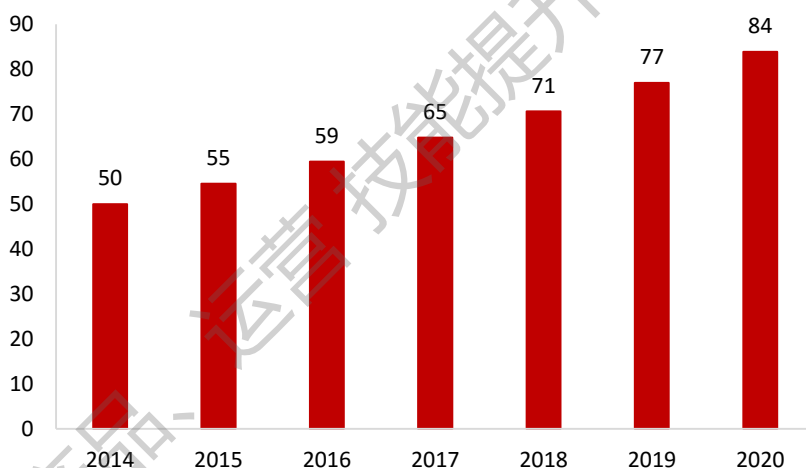
深度学习 FPGA 可以不再依赖于冯·诺依曼架构，而能够利用分布式片上存储器以及深度流水线并行，完美地契合了深度学习大计算量的要求；同时 FPGA 支持部分动态重新配置，这一特性大大降低大规模深度学习存储读取数据的成本；在算法层面 FPGA 给深度学习开拓了另一种思路：GPU

等固定架构设计遵循软件执行模型，需要算法进行适应，但 FPGA 较少强调算法去适应某固定计算框架，从而给算法留下更大的自由空间和发挥余地。

然而，FPGA 在展现架构优势的同时也存在不小的弊端，首先就是 FPGA 对于算法的要求更加宽泛，要求研究人员花费大量的时间去编译和完善；同时，FPGA 的硬件编辑语言十分复杂，这会影响 FPGA 应用于深度学习过程中的效率。

FPGA 高性能、低能耗以及可硬件编程的特点使其适用范围得以扩大，目前 FPGA 主要应用于通讯、医疗电子、安全、视频、工业自动化等领域。广阔的应用范围也拉动着 FPGA 未来庞大的市场规模。据 Gartner 统计，2014 年全球 FPGA 市场规模达到 50 亿美元，2015-2020 年的年均复合增长率为 9%，到 2020 年将达到 84 亿美元。

图 26：全球 FPGA 市场规模保持较快增长（单位：亿美元）

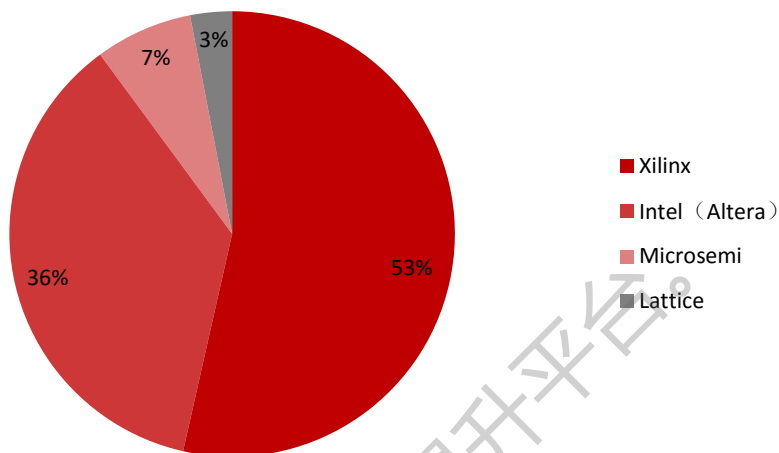


数据来源：Gartner，东方证券研究所

### 3.2 双寡头垄断 FPGA 市场

面对 FPGA 巨大的增长潜力，国际巨头纷纷尝试进入这一市场，据统计全球共有 60 多家公司先后出资数十亿美元，试图在 FPGA 行业占领一席之地，但目前全球 FPGA 市场主要被 Altera 和 Xilinx 瓜分，两家公司合计占有近 90% 的市场份额，合计专利达到 6000 多项，剩余份额被 Lattice 和 Microsemi 两家占据，合计共有超过 3000 项专利。技术专利的限制和漫长的开发周期使得 FPGA 行业形成了很高的壁垒，这也进一步巩固了 Altera 和 Xilinx 两家公司的优势地位和盈利水平。

图 27：2016 年 FPGA 市场份额分布



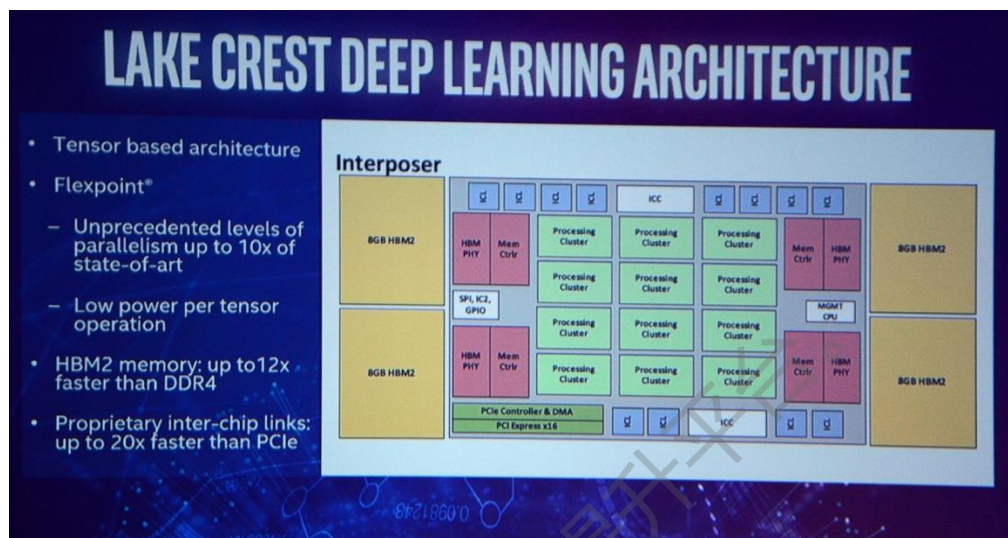
数据来源：Paul Dillien，东方证券研究所

Xilinx 和 Altera 公司在深度学习方面都取得了丰硕成果。Xilinx 提出未来深度学习处理器一定要经过模型压缩、模型定点化和编译三大步骤，并且拥有针对神经网络的专用结构。模型压缩使用户可以尽可能使用片上存储来存储深度学习算法模型，减少内存读取，以此大幅度减少能耗；模型定点化能够减少乘法器的大小；编译能够针对开发人员的具体要求进行特殊化处理从而实现更加高效的计算。Xilinx 提出的深度学习趋势为 FPGA 在这一领域的广泛应用打下坚实的理论基础。

2016 年初英特尔宣布以 167 亿美元的高价宣布收购 Altera 公司。英特尔作为在数据处理市场占据超过 95% 市场份额的巨头，一直在相关业务领域寻找新的增长点。目前的收购行为无疑表明英特尔将推动 FPGA 与 CPU 的整合，在未来的深度学习领域利用 FPGA 的硬件可编程性，在工作负载和计算需求发生波动时通过改变算法提高计算速度，同时维持较低功耗。在去年 11 月，英特尔发布了一款叫做 Nervana 的 AI 处理器，这个项目代码为“Lake Crest”，将会用到 Nervana Engine 和 Neon DNN 相关软件，这款芯片可以加速各类神经网络，例如谷歌 TensorFlow 框架，芯片由所谓的“处理集群”阵列构成，相对于浮点运算，这种方法所需的数据量更少，因此带来了 10 倍的性能提升。

图 28：英特尔 Lake Crest 架构





数据来源：Intel，东方证券研究所

从两家 FPGA 巨头的动作可以看出，由于 FPGA 在计算能力和灵活性上大大弥补了 CPU 的短板，从而未来在深度学习领域 CPU+FPGA 的组合将成为重要的发展方向。

### 3.3. 国内 FPGA 产业孜孜求索

2014 年中国 FPGA 市场规模已经达到 15 亿美元，占全球市场份额的三分之一，中国作为全球最大的通讯和军工市场之一，为了满足经济发展和通讯，尤其是国防等的需要，预计中国未来的 FPGA 市场需求量还会继续扩张。

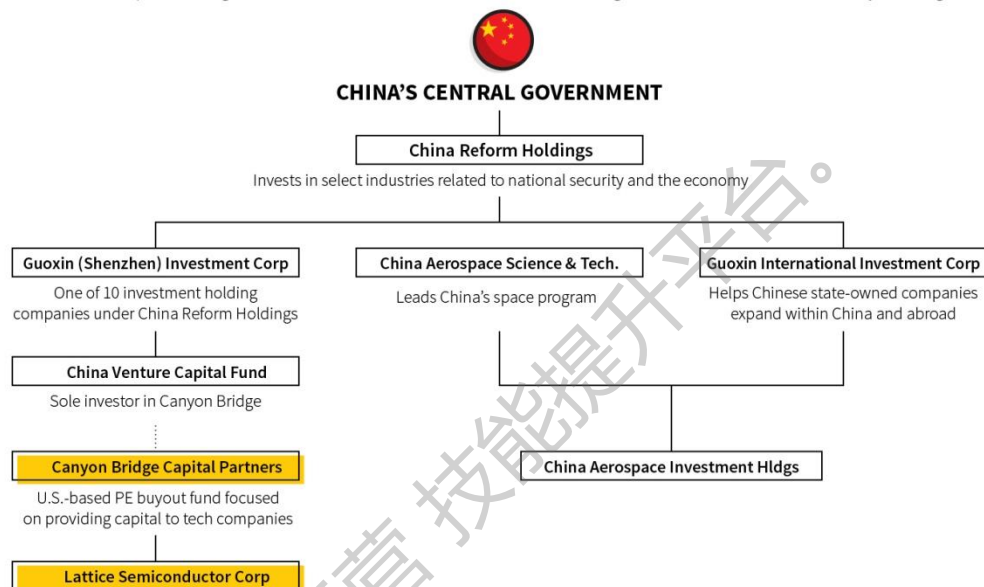
虽然政府多年来在此领域投入数百亿的科研经费，但由于美国对于技术专利的限制和 FPGA 高耸的技术门槛，国内 FPGA 探索的进程十分艰难，在产品性能、功耗、容量和应用领域上都存在较大差距。目前国内较为知名的 FPGA 相关公司仅有同创国芯、京微雅格、高云等。

2016 年 11 月，美国莱迪思半导体（Lattice）宣布，将被 Canyon Bridge Capital Partners 收购。后者是一家新成立的私募股权公司，唯一的投资人是 China Venture Capital Fund 的一家子公司，China Venture Capital Fund 则隶属于中国国新基金，这笔交易规模达 13 亿美元，若交易顺利达成，有望帮助国内企业在 FPGA 领域实现弯道超车。

图 29：Canyon Bridge Capital Partners 拟收购 Lattice

## Follow the money

Canyon Bridge Capital Partners, a buyout fund that agreed to buy U.S.-based chip maker Lattice Semiconductor for \$1.3 billion, is funded by cash tied to China's central government and also has indirect links to its space program, Chinese corporate filings show. The chart below shows how the Chinese government is connected to Canyon Bridge.



Sources: State Administration for Industry and Commerce of the People's Republic of China; U.S. Securities and Exchange Commission; China Reform Holdings promotional material.

C. Chan, 22/11/2016

REUTERS

数据来源：Business Insider、东方证券研究所

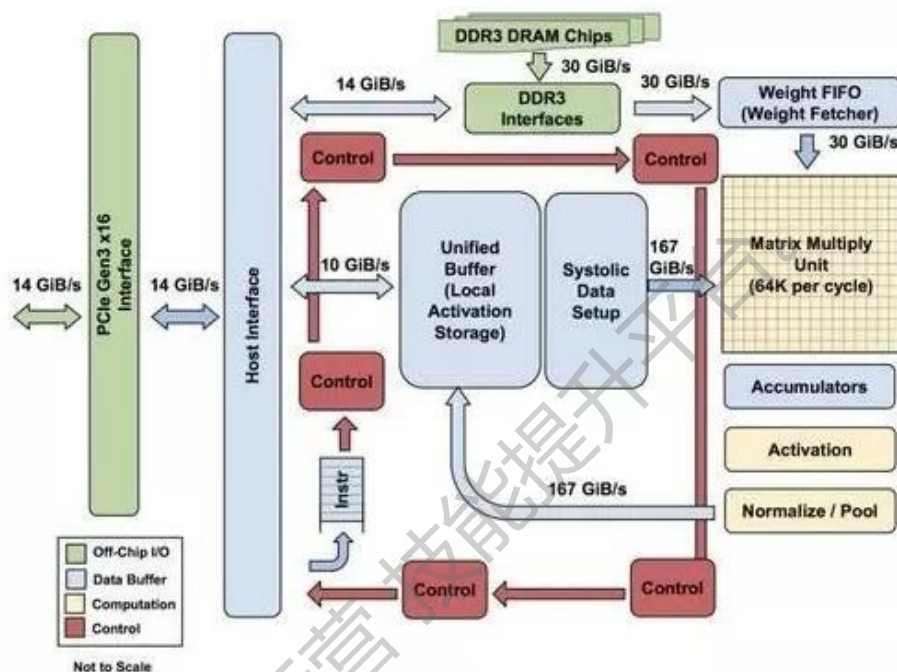
## 四、ASIC：有望成为主流趋势

为深度学习量身定制的 ASIC 芯片将在计算速度和功耗上赶超 GPU 和 FPGA，并随着人工智能渗透率的不断提升，未来在智能手机、物联网、车联网等领域，人工智能芯片将得到广泛应用，广阔的市场空间为 ASIC 大规模量产创造了可能。随着大规模量产条件下单片成本大幅下降，ASIC 可能会成为未来深度学习领域的主流芯片。

目前，科技巨头纷纷在 ASIC 深度学习芯片上发力，随着 AlphaGo 横扫人类顶尖棋手，谷歌在 AlphaGo 中应用的 ASIC 产品 TPU (Tensor Processing Unit) 最为受到业界的热捧，谷歌于 2016 年 Google I/O 大会上正式介绍第一代 TPU 产品，并于今年 4 月首次发表论文披露了 TPU 的详细架构和技术细节，根据谷歌的统计，CPU+TPU 的方案比 CPU+GPU 方案提高单位能耗计算能力 30~80 倍，提高计算速度 15~30 倍，适用于 Google 平台上 95% 的神经网络应用场景。

在今年 5 月的开发者 I/O 大会上，Google 正式公布了第二代 TPU，又称为 Cloud TPU，其最大的特色在于相比初代 TPU，它既可以用于训练神经网络，又可以用于推理，这既为推理阶段进行了优化，也为训练阶段进行了优化。在性能方面，第二代 TPU 可以达到 180TFLOPs 的浮点性能，和传统的 GPU 相比提升 15 倍，更是 CPU 浮点性能的 30 倍。

图 30：谷歌 TPU 内部架构



数据来源：Google，东方证券研究所

苹果公司正在研发一款名为“苹果神经引擎(Apple Neural Engine)”的 AI 专用芯片，该芯片定位于在本地设备上处理 AI 任务，旨在将主处理器和图像处理器巨大的计算量分开，把面部识别、语音识别等 AI 相关的任务卸载到 AI 专用模块上处理，以提升 AI 算法效率，并延长电池寿命，未来其可能应用于自动驾驶、Siri 语音助手及增强现实(AR)技术领域，未来还有可能嵌入 iPhone、iPad 等设备中，该芯片已在原型机中进行了测试。苹果有望在今年六月即将召开的年度开发者大会上公布 AI 芯片的研发进展。

国内在深度学习 ASIC 领域也不断取得突破进展。北京中科寒武纪科技有限公司研发了国际首个深度学习专用处理器芯片(NPU)，NPU 采用了“数据驱动并行计算”的架构，特别擅长处理视频、图像类的海量多媒体数据，其具有类似 GPU 的并行计算特点，但相比于 CPU，NPU 可以在线性代数运算上有更高的效率，但功耗上面可以比 CPU 低很多。目前寒武纪芯片 IP 指令集已扩大范围授权集成到手机、安防、可穿戴设备等终端芯片中，2016 年就已拿到一亿元订单。

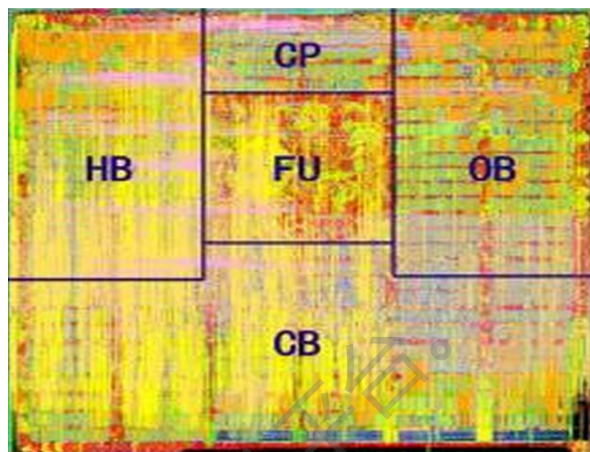
图 31：寒武纪芯片

图 32：寒武纪 2 号 DaDianNao 版图





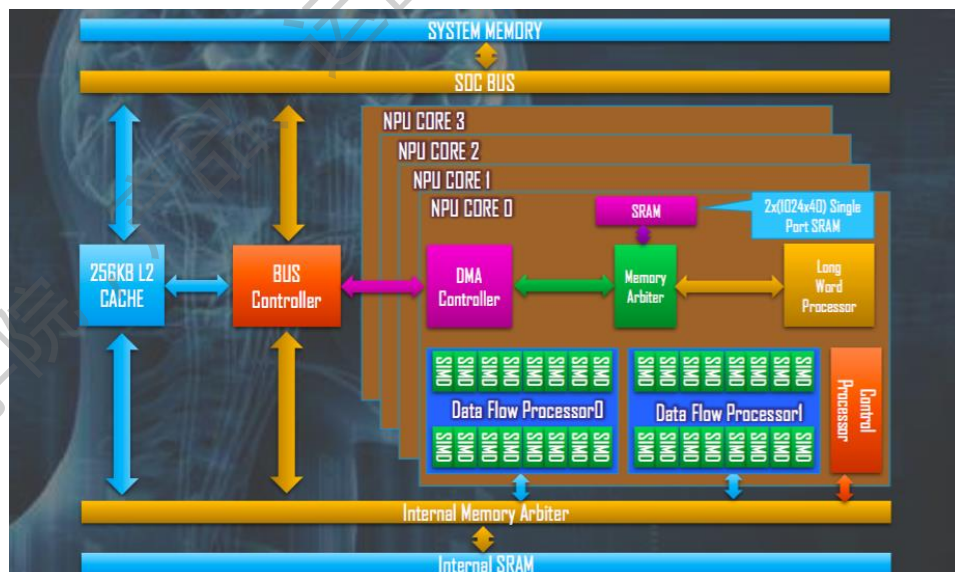
数据来源：雷锋网，东方证券研究所



数据来源：雷锋网，东方证券研究所

2016 年，中星微也推出了量产的 NPU 芯片“星光智能一号”，出货量主要集中在安防摄像领域，其中包含授权给其他安防摄像厂商部分，未来将主要向车载摄像头、无人机航拍、机器人和工业摄像机方面进行推广和应用。

图 33：中星微 NPU 架构图



数据来源：中星微，东方证券研究所

## 五、类脑芯片：超越“冯·诺依曼”架构的新思路

类脑芯片是一种基于神经形态工程、借鉴人脑信息处理方式、旨在打破“冯·诺依曼”架构束缚，适于实时处理非结构化信息、具有学习能力的超低功耗新型计算芯片。可以说类脑芯片是更加接近人工智能目标的芯片，其力图在基本架构上模仿人脑的工作原理，使用神经元和突触的方式替代传统“冯·诺依曼”架构体系，使芯片能够进行异步、并行、低速和分布式处理信息数据的能力，同

时具备自主感知、识别和学习的能力。

类脑芯片将主要实现两大突破，一是突破传统“执行程序”计算范式的局限，有望形成“自主认知”的新范式；二是突破传统计算机体系结构限制，实现数据并行传送、分布式处理，能够以极低的功耗实时处理海量数据。

类脑芯片实时海量数据处理及极低能耗的特性预示着其广阔的市场前景。根据 Markets and Markets 推测，如果类脑芯片能够顺利进入消费级应用，到 2022 年其市场规模将达到千亿级美元水平，消费终端将占整体市场的 98.17%，其他主流应用包括国防安全、工业自动化、航空航天等领域。

图 34：2022 年类脑芯片不同类型终端应用占比



正是由于类脑芯片巨大的发展潜力和广阔的市场前景，各国政府及科技巨头都在大力推动类脑芯片的研发进程，包括美国、日本、德国、英国、瑞士等发达国家已经制定相应的发展战略，中国的类脑科学研究项目目前也已经正式启动。

图 35：各国类脑计算研究项目列表

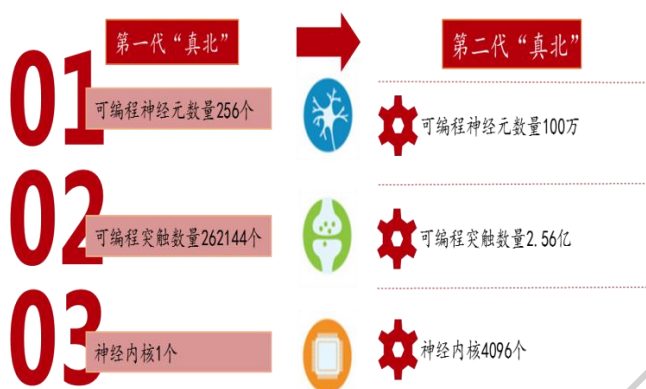
时间	国家	机构	项目名称	开展原因	进展
2003年	日本	日本政府	脑科学与教育	将脑科学研究作为国家教育发展的一项战略任务	正在面向教育理论和实际应用进行研究
2012年	美国	美国国防高级研究计划局	DARPA类脑图像处理项目	在情报、监视与侦查数据中，图像与视频占据很大比重，而传统的图像处理器受器件和架构的制约，其性能虽然不断增长，但无法满足日益增长的战争画面和视频处理需要	与密歇根大学合作开发处理速度比目前图像处理器快1000倍，但功耗仅为万分之一的类脑图像处理器，密歇根大学将在四年内分两阶段完成
2013年	美国	美国国立卫生研究院	BRAIN计划	推动美国神经技术及脑科学研究	工作组计划在未来五年投资将达到每年4亿美元，随后5年为每年5亿美元
2013年	欧盟	欧盟未来技术项目	欧盟人脑计划	旨在建立一套基于神经科学的最新的、革命性的信息通信技术，建造一种模拟神经元功能的芯片，并将这种芯片用于建造超级计算机系统	该计划将持续十年，整体投资11.9亿欧元
2015年	美国	情报高级研究计划局	大脑皮层网络机器智能项目	试图通过数据科学与神经科学的结合，通过人类大脑逆向工程算法快速推进机器学习阿赫人工智能研究，以提高对复杂信息的处理能力	该项目计划执行期五年，分三个阶段完成，各阶段将会涉及人脑神经解剖学和神经生理学研究，以增进对基于感觉信息处理的大脑皮层计算能力的认识
2015年	中国	中国科技部	中国脑计划	从认识脑、保护脑和模拟脑三个方向全面启动。制定中国的脑科学和类脑研究方案	清华大学、中国科学院已经成立类脑研究中心
2017年	中国	中科院、复旦大学、百度、微软等	类脑智能技术及应用国家工程实验室	建立脑认知和脑模拟技术研究与实验平台	实验室成立





2014 年 8 月，IBM 公司推出第二代 TrueNorth 芯片，采用 28nm 硅工艺技术，包括 54 亿个晶体管和 4096 个处理核，相当于 100 万个可编程神经元，以及 2.56 亿个可编程突触。TrueNorth 每个处理核中包含约 120 万个晶体管，大多数晶体管用作数据存储、以及与其他核的通信，因此芯片的工作方式类似于人脑的神经元和突触之间的协同。与一代相比，二代 TrueNorth 芯片性能大幅提高，且处理核体积仅为第一代的 1/15。目前，IBM 公司已经利用 16 颗 TrueNorth 芯片开发出一台神经元计算机原型，具有实时视频处理能力。

图 38：第一代 IBM TrueNorth 芯片与第二代比较



数据来源：东方证券研究所

图 39：IBM 神经元计算机包含 16 颗 TrueNorth 芯片



数据来源：互联网，东方证券研究所

## 六、人工智能芯片在云端与终端携手共进

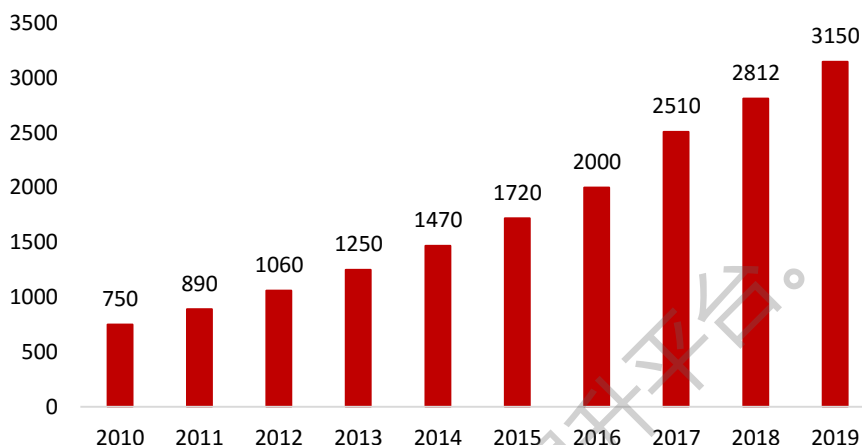
在全球智能化发展的浪潮中，人工智能已经成为未来发展的重要领域，根据经济学人的相关调查，在 2015 年以后，市场对于人工智能领域的关注度呈现指数级增长，人工智能技术有望引领新一轮科技革命，对于各个科技强国及科技巨头，如何构建最佳的架构和系统来处理 AI 工作所必需的海量数据是重中之重，从最初的 CPU 到目前应用较为广泛的加速计算 GPU、FPGA，再到前沿的 ASIC、类脑芯片，芯片作为人工智能技术的核心技术环节，决定了整个领域未来的发展方向。

### 6.1. 云端 AI 芯片领域百家争鸣

人工智能技术的发展跟数据量的飞跃式发展有密不可分的关系。根据 IDC 报告显示，预计到 2020 年全球数据总量将超过 40ZB，而这一数据量是 2011 年的 22 倍。并且在过去的几年，全球的数据量以每年 58% 的速度增长。面对如此庞大的数据量，目前平均每年仅有 0.4% 的数据得到了良好的分析利用，因此，进一步发展人工智能关键之一就是增强数据挖掘的“纵深”，分析更深层面、更大规模的数据。

另根据 IDC 的调查，随着云计算技术的不断发展，和其超强计算、成本较低等特性被大众所挖掘，58% 的受调查企业计划使用基于网络的云计算服务，而这一比例远超 2014 年的 24%。云计算的市场规模也在逐渐扩大，据 Gartner 的统计，到 2019 年，全球云计算市场规模将达到 3150 亿美元，远超当前的 1720 亿美元。

图 40：全球云计算市场规模（亿美元）



数据来源：Gartner，东方证券研究所

因此，人工智能关键技术是在云计算和大数据日益成熟的背景下取得了突破性进展，云计算为人工智能提供平台，而大数据为人工智能提供信息来源。目前各大科技巨头看好未来人工智能走向云端的发展态势，纷纷在自有云平台基础上搭载人工智能系统，以期利用沉淀在云端的大数据挖掘价值。

图 41：云计算平台人工智能功能

公司	云计算平台	使用芯片	人工智能功能介绍
IBM	Watson	POWER7、8	IBM将人工智能与商业分析功能融合在Watson平台上，使得Watson能够帮助用户进行问题分析、开发方案、监督学习、肿瘤治疗及临床试验匹配
微软	Azure	英特尔CPU+Nvidia GPU	微软将人工智能加载在Azure上，以期为物联网打造安全、可靠、灵活、高效的云后台。通过事件监控、数据存储、分析转化、结果呈现和知道决策等环节融会贯通
亚马逊	AWS	英特尔Broadwell-EX 至强处理器+GPU集群	亚马逊将在AWS上提供“亚马逊机器学习”服务，让AWS服务的网站变得更聪明，增加网络应用的预测和分析能力
谷歌	Google Cloud Platform	IBM POWER8处理器+TPU	人工智能已经成为谷歌云计算战略的核心，谷歌目前公布的软件功能包括可以提取文本内容含义且能将语音内容转化为文本，从而帮助用户分析网站信息，协助客服电话
阿里巴巴	阿里云	NVIDIA Tesla M40 GPU+Intel酷睿第5代CPU	基于阿里云开发了国内第一个人工智能平台“DTPAI”，通过此平台，开发者能够对用户行为、行业走势等进行预测，该平台已经集成阿里巴巴核心算法库
百度	百度开放云	英特尔至强E5-2600 v4系列处理器+Altera FPGA集群	百度开放云是堆栈为三层架构：最下面是云计算层，其上为大数据应用层，最上为人工智能层，大数据借助人工智能技术，将具体应用与云计算技术结合成为行业解决方案

数据来源：互联网，东方证券研究所

## IBM Watson

IBM Watson 由 IBM 历经 4 年时间研发，并于 2011 年参加美国电视问答节目《Jeopardy》勇夺第一而一战成名。目前 Watson 已经发展成为集分析、发现、诊断、教学、侦查等功能为一体的综合性认知计算系统。

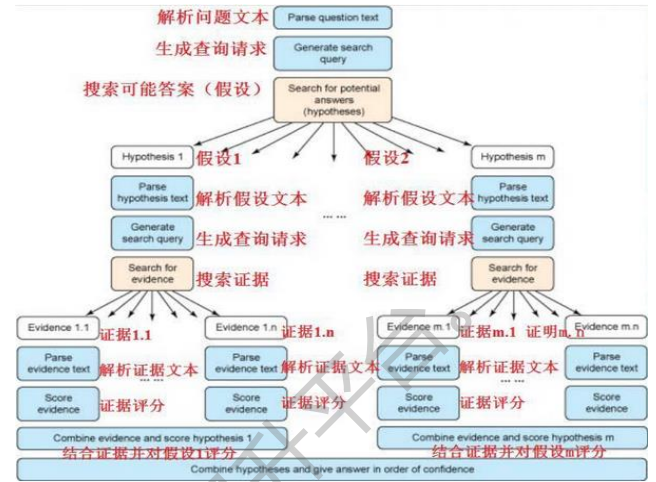
图 42：2011 年 Watson 参加节目《Jeopardy》并取得冠军

图 43：Watson 产生答案流程





数据来源：雷锋网，东方证券研究所



数据来源：雷锋网，东方证券研究所

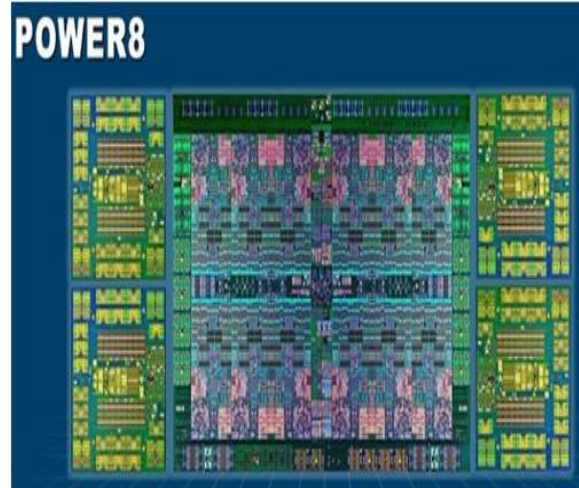
Watson 惊艳的计算、分析能力来源于其背后强大的处理器，Watson 系统搭载 IBM 的 POWER7 处理器，采用 45nm 生产工艺，提供 4-8 核型号，每个核心支持 4 线程，即每颗处理器最高同时支持 32 线程运算。POWER7 处理器的创新之处在于其能够随数据量变化调节内核运行方式并设置不同的线程模式，从而最大化工作效率。2013 年 IBM 推出 POWER8 处理器，升级制造工艺至 22nm，核心数量达到 12 个，无论在多线程能力、缓存容量，还是内存带宽上都有飞跃式的发展。

图 44：IBM POWER 处理器发展路径

图 45：POWER8 架构图

	POWER5 2004	POWER6 2007	POWER7 2010	POWER7+ 2012	POWER8
Technology	130nm SOI	65nm SOI	45nm SOI eDRAM	32nm SOI eDRAM	22nm SOI eDRAM
Compute					
Cores	2	2	8	8	12
Threads	SMT2	SMT2	SMT4	SMT4	SMT8
Caching					
On-chip	1.9MB	8MB	2 + 32MB	2 + 80MB	6 + 96MB
Off-chip	36MB	32MB	None	None	128MB
Bandwidth					
Sust. Mem.	15GB/s	30GB/s	100GB/s	100GB/s	230GB/s
Peak I/O	3GB/s	10GB/s	20GB/s	20GB/s	48GB/s

数据来源：IBM，东方证券研究所



数据来源：IBM，东方证券研究所

随着云计算的不断发展，越来越多的科技巨头选择将自己的服务器“云端化”。目前 IBM 公司正致力于开发 Watson 云端服务，旨在将 Watson 强大的理解、推理和学习能力在云平台上进行整合，由此扩大适用范围，增加客户数量。

## 微软 Azure

Azure 是微软公司于 2008 年发布的云计算操作系统，主要目标是为开发者提供平台以帮助其开发在云服务器、数据中心、Web 和 PC 上的应用程序。Azure 使用英特尔 CPU 作为其处理器，并于 2015 年 10 月宣布将通过 Azure 平台向全球客户提供基于 Nvidia 的 GPU 专业图形应用和加速计算功能，即在其处理器上加载 Tesla K80 GPU 加速器，旨在向客户提供高级计算级性能，以便满足更严苛的数据中心与高性能计算应用需求。

图 46：微软 Azure 功能



数据来源：微软，东方证券研究所

## 亚马逊 AWS

亚马逊的 Amazon Web Services (AWS) 于 2006 年推出，旨在以云端服务形式向企业提供 IT 基础设施服务。其主要优势是能够利用云端 IT 设施帮助企业节省大量前期资本基础设施费用。在云计算领域，亚马逊无疑是同行业中的佼佼者，根据 Synergy Research Group 于 2014 年发布的全球云计算调查报告，亚马逊 AWS 以 24% 的市场份额占比高居榜首，远超第二名微软 Azure 10% 的市场份额占比。

图 47：2014 年亚马逊 AWS 市场份额占比遥遥领先

图 48：亚马逊 AWS 能够提供的服务



### Cloud Infrastructure Services

#### MARKET SHARE, 2014

28% Amazon Web Services

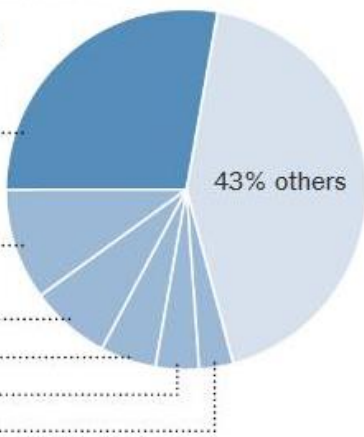
10 Microsoft

7 I.B.M.

5 Google

4 Salesforce

3 Rackspace



数据来源：Synergy Research Group，东方证券研究所



数据来源：Amazon，东方证券研究所

目前 AWS 使用英特尔最新 Broadwell-EX 至强处理器以提升平台分析能力及内存数据处理能力，并在 2013 年正式推出 GPU 实例服务。这项服务将使用 Nvidia Grid GPU 并行处理功能以支持计算量很大的应用。并且，AWS 内置 GPU 集群，提供 33.5 个计算单元，并搭载 2 块 Nvidia Tesla M2050 显卡，大幅度提升服务器的计算能力。

同时在 2014 年 4 月，亚马逊宣布将自主设计服务器芯片，并且完成对一家以色列半导体初创公司的收购，以期更好地满足服务器的计算需求。

### 谷歌云平台

谷歌在云计算市场姗姗来迟，但目前正在加快这部分业务的拓展力度。2015 年谷歌投入约 100 亿美元建设新的数据中心，据传谷歌的云计算将使用 IBM 的 POWER8 芯片，并将搭载在 AlphaGo 中使用过的 TPU 作为加速计算芯片，意在将云计算平台作为人工智能输出的渠道和平台。

图 49：谷歌云计算平台



Google Cloud Platform

数据来源：IT168，东方证券研究所

## 阿里云

阿里云创立于 2009 年，致力于为企业、政府等组织机构，提供最安全、可靠的计算和数据处理能力。同时，阿里云正在利用自身强大的数据支持和计算能力开拓更多的应用场景。2015 年 4 月中石化与阿里云正式建立技术合作关系，以期借助云计算和大数据的力量改变传统的石油化工业务；同时，阿里云致力于开展量子信息科学研究，并与英特尔、华大基因合作共建医疗应用平台，让我们看到了中国云计算未来的无限可能。

图 50：阿里云适用场景



数据来源：阿里云，东方证券研究所

阿里云使用芯片为英特尔 Xeon 系列芯片。采用 14nm 制造工艺，最大核心数量能够达到 22 个，并拥有最多 44 个逻辑线程，在计算性能和安全性上都有较大提升。并且在 2016 年 3 月，阿里云发布面向深度学习、3D 图形图像渲染及科学计算的新一代 HPC 平台。新平台采用 NVIDIA Tesla M40 GPU 超大规模加速器和 Intel 酷睿第 5 代 CPU，整体性能大幅飞跃。

图 51：阿里云新一代 HPC



数据来源：阿里云，东方证券研究所

## 百度开放云

百度开放云平台是百度 2015 年开发的基于大数据为企业、政府及个人提供云服务的云计算平台。使用处理器为英特尔 Broadwell-EP 架构的至强 E5-2600 v4 系列处理器，跟亚马逊所用处理器相同。虽然百度开放云较阿里云晚六年才进入市场，但是其平台优势还是十分明显：在硬件方面，从两三年前开始，百度就已经在用 GPU 代替 CPU 进行计算，百度很清楚的认识到了未来全球云计算的发展趋势，同时与 Altera 公司合作在云数据中心使用 FPGA 集群，百度作为最早一批在服务器中使用 FPGA 集群的科技巨头，满足了数据中心环境对高性能和灵活性的要求。

在功能方面，后入场的百度也精准得找到了云计算的发展趋势，即“云计算+大数据+人工智能”。搜索业务出身的百度在大数据上有明显的优势，且在大数据能力基础上进行人工智能布局也成为目前几乎所有云平台共同的发展方向。

图 52：百度与 Altera 合作建立 FPGA 集群



数据来源：电子产品世界，东方证券研究所

图 53：百度开放云功能



数据来源：百度，东方证券研究所

随着科技巨头服务器业务的蓬勃发展，我们已经能够看到一条清晰的发展趋势，即“云端化”+“AI 芯片集群化”。目前市场主流服务器提供商均采用“云平台”方式，以最高效的方式打包为用户提供最大化便利。同时，各大云计算平台已经开始搭载人工智能芯片，或建立人工智能芯片集群，以满足激增的数据量要求并提升平台计算能力。

## 6.2. 终端 AI 芯片领域初露头角

当前各大科技巨头在人工智能芯片领域的布局大多集中在云端 AI 芯片领域，在云端处理与 AI 相关的任务，虽然这种方法可以获得优异的处理性能，但在安全性和隐私性方面有所牺牲。对于广阔的消费电子市场，终端 AI 芯片领域未来有望放量。

从研发周期、生产成本、产品性能、功耗水平等多个角度进行衡量，最容易应用到终端消费领域的加速学习芯片是 GPU。谷歌高级研究员在预测深度学习未来趋势时谈到未来深度学习算法将会更高效，即使没有额外的硬件支持或是过高的内存开销，也能够廉价的移动设备上运行。消费级核心芯片的要求主要有两个：首先是体积和性能上能够满足移动智能设备的要求；其次是成本保持在较低水平。GPU 基于其较低的生产成本、较为成熟的生产工艺、强大的处理并行计算的能力及在

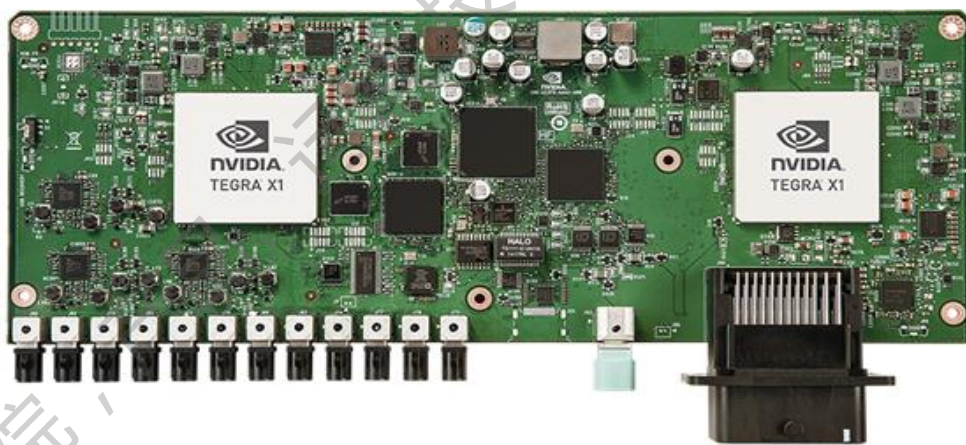
图像处理方面在智能手机端积累了较多应用基础，有望成为最可能抢滩消费级加速计算市场的核心芯片。而在性能要求较高的深度学习领域，例如智能机器人及类脑计算机的研发，FPGA、ASIC及目前正处于尝试阶段的类脑芯片的技术发展都有望对机器学习领域的进步做出巨大贡献。目前在汽车、机器人、家居等终端领域，人工智能技术已经开始得到应用，部分科技巨头也开始切入相应市场进行布局。

## 智能汽车

ADAS 系统是用机器部分甚至完全替代驾驶员的技术，ADAS 采用摄像头、激光雷达、毫米波雷达等多种感知手段，深度学习是自动驾驶系统的主要算法。ADAS 汽车需要强大的人工智能芯片，GPU 是当前的主流技术路线。

2015 年 5 月，Nvidia 发布的应用于无人驾驶汽车的车载计算平台 Drive PX，集成了两个 Tegra X1 处理器，该处理器采用 20nm 制造工艺，并集成两颗 GPU 和八颗 CPU 芯片，利用 GPU 强大的并行计算能力和图像处理能力进行汽车周边环境的探测和监控，实现环视系统、碰撞规避系统、行人检测系统和驾驶员状态监测系统等功能。

图 54：Nvidia Drive PX 车载计算平台



数据来源：百度，东方证券研究所

2016 年 6 月，Nvidia 又推出全新的 Drive PX2，采用 12 核处理器并搭载 2 颗基于 Pascal 架构的新一代 GPU——Tegra K1。该芯片除了强大的计算能力外，还配有深度学习功能，旨在最大化发挥 GPU 芯片的加速计算功能。

图 55：Nvidia Drive PX2 平台

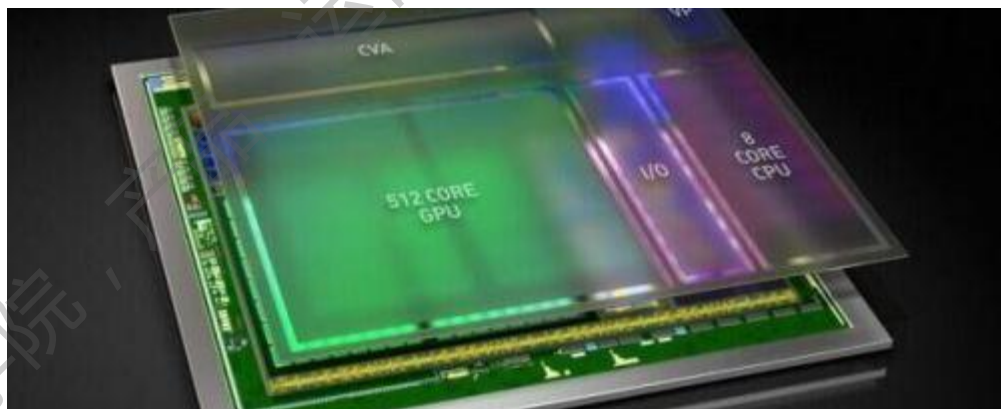




数据来源：Pconline，东方证券研究所

2016 年 9 月，Nvidia 进一步发布了一款针对自动驾驶技术和汽车产品的芯片 Xavier，Xavier 采用自定义的八核 CPU 架构，同时内建 Nvidia 全新 Volta GPU 架构，是自动驾驶汽车的计算机视觉加速器。Xavier 采用 16nm FinFET 工艺，在提升性能的同时降低功耗，Xavier 运算性能达到 20TOPS，功耗则只需 20 瓦。

图 56：Nvidia Xavier 芯片



数据来源：NVIDIA，东方证券研究所

面对 Nvidia 在汽车市场的不断突破，其他传统半导体巨头纷纷在智能汽车芯片领域发力。高通作为移动终端处理器的优势企业，于 2014 年 6 月发布智能汽车芯片——骁龙 602A。该处理器同样选择 CPU+GPU 模式，在快速处理数据信息的同时，提升地图的渲染效果，并降低处理器能耗。同时高通还推出 Snapdragon 820 平台，以单芯片整合 64 位 CPU、GPU 和 DSP，实现 ADAS 的大部分功能。

图 57：高通发布智能汽车芯片 602A

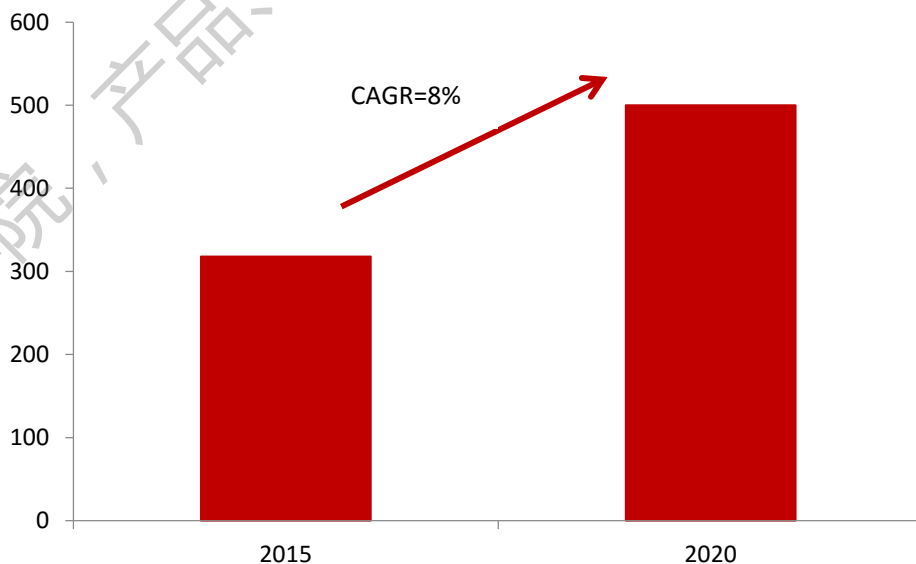


数据来源：IT168，东方证券研究所

其他厂商如恩智浦、意法半导体、三星、英特尔等均发挥自身优势，推出各种智能汽车应用方案，以期在智能汽车时代找到新的业绩增长点。

国内车载芯片市场规模同样巨大，2015 年国内汽车电子芯片市场规模为 318 亿元，随着政府大力支持国内厂商自主研发芯片，获取产业链上高附加值，未来自主研发汽车芯片企业有望实现突破，打入国际主流厂商供应链，逐步取代进口芯片。我国汽车电子芯片行业需求将高速增长，预计 2021 年，我国汽车电子芯片行业需求规模将达到约 500 亿元。

图 58：国内汽车电子芯片市场规模



数据来源：中投顾问产业研究中心，东方证券研究所

虽然中国智能汽车市场起步较晚，但是已经有相关厂商在该领域做出重大突破。地平线机器人公司于今年 3 月推出面向自动驾驶的“雨果”平台，在硬件方面，地平线计划将 NPU 集成到平台之上，

预计计算性能将比目前提升 2-3 个数量级。

智能汽车这一潜力无限的“蓝海”市场有极大的发展空间，而芯片作为整个智能汽车行业的先导将成为众多半导体巨头竞相开拓的高地。目前智能汽车芯片均采用 CPU+GPU 形式，以 CPU 强调逻辑处理，并利用 GPU 强大的并行计算和图像处理能力。国内的地平线已经开始尝试在智能驾驶平台搭载 ASIC 芯片。定制芯片无疑能将数据处理速度提升到一个新的高度，并将能耗维持在相对较低水平，但鉴于其研发周期长且成本高昂，我们认为未来几年 CPU+GPU 仍然是智能汽车芯片的主流方案。

## 智能机器人

伴随着信息技术的不断进步，人类已经在越来越多的场景中试图用机器人进行工作。根据 Markets and Markets 的统计，到 2020 年智能机器人的全球市场规模将达到 78.5 亿美元，年复合增长率为 19.22%。

在工业机器人领域，CNC（数控机床）对机器视觉和机器手臂的大范围应用催生了机器人内部控制元件的蓬勃发展。飞思卡尔已经推出融合 GPU 的 Vybrid 处理器，以替代原有的 DSP 芯片，以为开发商提供更加精准的机器视觉处理方案。另一方面，出于工厂产线精密化，机器手臂多轴化的发展趋势，赛灵思已经开始布局将 FPGA 芯片应用于机器手臂，利用 FPGA 芯片的扩展性和强性能应对设计日益复杂的多轴马达运算需求。相较于单节点 DSP 芯片，FPGA 不仅能将启动电流环耗时由平均 200 微秒降低至 50 微秒，大幅提高马达运作速度，同时能够实现单一颗芯片控制多轴机器臂的任务。

图 59：飞思卡尔 Vybrid 处理器



数据来源：互联网，东方证券研究所

图 60：赛灵思 FPGA 芯片



数据来源：互联网，东方证券研究所

同时，已经有科技公司在消费级机器人市场发力。今年 5 月底日本夏普公司开始发售一款机器人手机 RoBoHoN，这款颠覆传统设计观念的智能手机，不仅在外形上是一个机器人，还具备直立行走、跳舞、语音互动、人脸识别及投影等功能，同时兼顾手机接打电话，收发邮件等功能。

图 61：夏普机器人手机 RoBoHoN





数据来源：IT 之家，东方证券研究所

RoBoHoN 采用高通骁龙 400 处理器，该处理器芯片搭载八个 CPU 和 Adreno 305 GPU，能将图形处理性能提高 50%。从处理器的性能和夏普公布的信息来看，这款机器人手机并非想要在性能上与其他产品一决高下，而是意图在智能手机高度同质化的当下，利用机器人概念开辟细分市场；同时，RoBoHoN 为消费级机器人产品开拓了崭新而庞大的应用领域，也为人工智能的产品落地实现了重要突破。

## 智能音箱

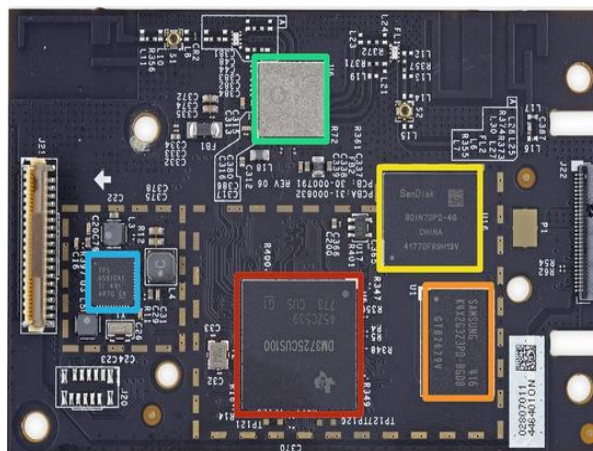
物联网应用的蓬勃发展将智能家居推到了历史前台。智能家居作为推进物联网发展至关重要的一环，蕴藏着巨大的发展空间，根据物联中国的测算，到 2018 年，全球智能家居的市场规模将达到 710 亿美元；而中国近几年在智能家居领域将呈现爆发式增长态势，平均年增速在 50% 左右，2018 年市场规模将达到 225 亿美元，占全球智能家居市场份额的 31%。

目前，亚马逊、谷歌、苹果、三星、微软等科技巨头纷纷开始布局智能家居市场。其中亚马逊推出智能音箱 Echo，在支持音箱功能的同时，更支持语音搜索、购物、提醒等多项操作，在未来，Echo 还将被加入更多与其他智能家居的兼容性，使 Echo 能够通过用户的指令完成对特定家居的操作。Echo 音箱主要芯片包括德州仪器的 DSP 和集成电源管理 IC，三星的 RAM，SanDisk 的 4GB 闪存和高通的 Wi-Fi、蓝牙模块。

图 62：亚马逊 Echo 音箱基本构造



图 63：Echo 音箱主板芯片构成





数据来源：物联中国，东方证券研究所

数据来源：爱板网，东方证券研究所

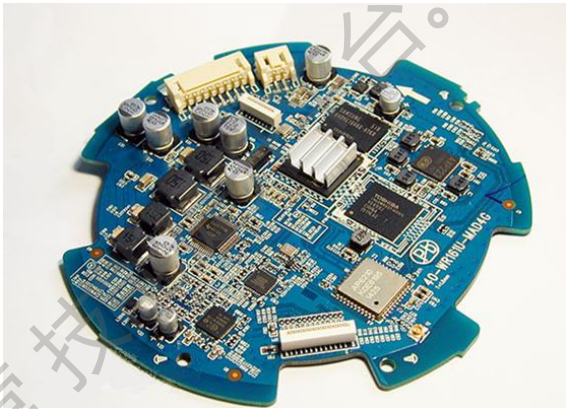
国内科技巨头在该领域也一定有所突破。京东与国内最大语音技术公司科大讯飞联合开发叮咚音箱，能够在为用户提供音箱功能的同时，支持语音控制，并致力于在未来成为智能家居的集中控制中心。音箱主芯片采用全志四核 Cortex-A7 CPU，并内置 Mali400 GPU，旨在发挥其计算及音频处理功能。

图 64：京东&科大讯飞叮咚音箱



数据来源：互联网，东方证券研究所

图 65：叮咚音箱主板构造



数据来源：搜狐科技，东方证券研究所

随着人工智能技术的逐渐成熟，应用场景和落地产品将日益丰富，许多传统行业有望乘着人工智能技术的东风开拓新的业务模式，挖掘新的市场机会。芯片作为人工智能技术的核心环节，云端 AI 芯片与终端 AI 芯片将双轮并举，进一步加速各行业智能化渗透进程。

图 66：人工智能芯片及应用

行业	使用领域	产品名称	半导体公司	芯片类型	功能
服务器	认知计算平台	Watson	IBM	POWER7、8处理器	POWER7采用45nm生产工艺，每核心支持4线程，即每颗处理器最高同时支持32线程运算；POWER8制造工艺至22nm，核心数量达到12个
	云计算	Azure	微软	英特尔CPU+Nvidia GPU	Azure使用英特尔CPU作为其处理器，并在其处理器上加载Nvidia Tesla K80 GPU加速器，旨在向客户提供高级计算级性能
	云计算	AWS	亚马逊	英特尔Broadwell-EX 至强处理器+GPU集群	AWS内置GPU集群，提供33.5个计算单元，并搭载2块Nvidia Tesla M2050显卡，大幅度提升服务器的计算能力
	云计算	Google Cloud Platform	谷歌	IBM POWER8+TPU	人工智能已经成为谷歌云计算战略的核心，功能包括可以提取文本内容含义且能将语音内容转化为文本，从而帮助用户分析网站信息，协助客服电话
	云计算	阿里云	阿里巴巴	NVIDIA Tesla M40 GPU+Intel酷睿第5代CPU	阿里云发布面向深度学习、3D图形图像渲染及科学计算新一代HPC平台，采用NVIDIA Tesla M40 GPU加速器和Intel酷睿i5 CPU，整体性能大幅飞跃
	云计算	百度开放云	百度	英特尔至强E5-2600 v4系列处理器+Altera FPGA集群	百度已经用GPU代替CPU进行计算，并与Altera公司合作使用FPGA集群，极大满足数据中心环境对高性能和灵活性的要求
智能汽车	无人驾驶	车载计算平台Drive PX	Nvidia	Tegra X1处理器	采用20nm制造工艺，集成两颗GPU和八个CPU芯片，利用GPU并行计算和图像处理能力进行汽车周边环境的探测和监控
	无人驾驶	车载计算平台Drive PX2	Nvidia	Tegra K1处理器	配有深度学习功能，最大化发挥GPU芯片在加速计算上的能力
	智能汽车	智能汽车芯片	高通	骁龙602A	选择CPU+GPU模式，在快速处理数据信息的同时，提升地图的渲染效果，并降低处理器能耗
	ADAS	Snapdragon 820平台	高通	骁龙820	以单芯片整合64位CPU、GPU和DSP，实现ADAS的大部分功能
	V2X	V2X芯片组	恩智浦	RoadLINK芯片	恩智浦与德尔福汽车公司展开合作，并顺利量产V2X通信芯片
	V2X	V2X芯片组	意法半导体	V2X芯片	意法半导体与Autotalks合作开发V2X芯片组，并计划于2017年完成第二代V2X芯片组的大规模部署
	无人驾驶	“雨果”平台	地平线机器人	NPU	地平线计划将NPU——一个专门为自动驾驶数据处理而研发的芯片，集成到平台之上，预计计算性能将比目前提升2-3个数量级。
智能机器人	工业机器人	Vybrid处理器	飞思卡尔	Vybrid处理器	Vybrid处理器上搭载GPU芯片，以替代原有的DSP芯片
	机器人手臂	FPGA芯片	赛灵思	FPGA芯片	利用FPGA芯片的扩展性和强性能应对设计日益复杂的多轴马达运算需求
	消费级机器人	机器人手机RoBoHoN	夏普	骁龙400处理器	搭载八核CPU和Adreno 305 GPU，能将图形处理性能提高50%
智能音箱	智能音箱	Echo音箱	亚马逊	德州仪器DSP	Echo音箱主要芯片包括德州仪器的DSP和集成电源管理IC，三星的RAM，SanDisk的4GB闪存和高通的Wi-Fi、蓝牙模块
	智能音箱	叮咚音箱	京东&科大讯飞	全志四核Cortex-A7 CPU	音箱主芯片采用全志四核Cortex-A7 CPU，并内置Mali400 GPU，旨在发挥其计算及音频处理功能

数据来源：互联网，东方证券研究所

### 投资建议

虽然我国在人工智能领域的积淀时间相对较短，但发展迅速：国内领先企业已经在无人驾驶、车联网、机器人、智能家居和云端服务器领域取得了较多突破进展；同时中国人工智能市场份额年增速高达 50%，远超全球平均水平的 19.7%。相信随着相关知识产权的不断开放和技术的不断积累，未来我国在人工智能芯片领域的发展速度将大幅提高。

目前，国内已经有部分企业在沿人工智能产业链上进行布局：在核心芯片领域，建议关注中科曙光(603019，未评级)、全志科技(300458，未评级)、景嘉微(300474，未评级)、通富微电(002156，未评级)；在安防领域，建议关注海康威视(002415，未评级)、大华股份(002236，未评级)、富瀚微(300613，未评级)。

图 67：A 股上市公司切入人工智能领域情况

代码	公司名称	人工智能产业链	公司主营业务	人工智能发展情况
603019	中科曙光	基础芯片	研究、生产、制造高性能计算机、通用服务器及存储产品	背靠中科院，与国内首款神经网络处理器研发机构寒武纪展开合作；发布基于Nvidia GPU的深度学习平台
300458	全志科技	基础芯片	智能应用处理器SoC和智能模拟芯片设计厂商	依托芯片设计先发优势，募资开展车联网智能终端应用处理器芯片与模组研发及应用云建设项目，公司芯片产品也被用作叮咚音箱主芯片
300474	景嘉微	基础芯片	主要涉及图形显控及小型专用化雷达两大领域	公司研发的JM5400图形芯片打破国外芯片在我国军用GPU领域的垄断，实现了军用GPU国产化
002156	通富微电	基础芯片	国内封测产业三巨头之一	公司通过收购AMD旗下子公司，切入GPU封装和测试领域。
002415	海康威视	安防	国内视频监控领域领先企业	携手NVIDIA和Movidius发布了基于深度学习技术的从前端到后端全系列智能安防产品，推出“深眸”、“脸谱”、“神捕”等新产品
002236	大华股份	安防	国内视频监控领域领先企业	整合了上百台高性能计算机，建设深度学习计算集群，现已完成了深度学习计算中心的建设。
300613	富瀚微	安防	国内领先的视频监控芯片及解决方案厂商	公司14年底与赛蓝科技合作推出新一代WIFI智能可视门铃，基于深度学习算法，将人工智能适当介入安防视频监控，包括对目标分类、区分识别等不同场景。

数据来源：东方证券研究所整理

结合公司整体业务和人工智能芯片领域的状况，我们建议重点关注中科曙光、全志科技、景嘉微、通富微电、富瀚微等公司。

中科曙光(603019，未评级)：超级计算机龙头，牵手寒武纪进军深度学习领域。

- 1) 公司作为国内超级计算机龙头企业，在服务器市场占据优势，积极切入云计算领域，在全国建设多处云计算中心，为发展人工智能业务打下了良好基础。
- 2) 公司背靠中科院，与国内首款神经网络处理器研发机构寒武纪展开合作，已发布基于 Nvidia GPU 的深度学习平台。

全志科技(300458，未评级)：芯片设计领域厚积薄发，人工智能崭露头角。

- 1) 公司在芯片设计领域积累多年，在超高清视频编解码、CPU/GPU 多核整合、先进工艺高集成度等方面处于业界领先水平。
- 2) 公司依托芯片设计先发优势，募资开展车联网智能终端应用处理器芯片与模组研发及应用云建设项目，公司芯片产品也被用作叮咚音箱主芯片，在人工智能领域崭露头角。

景嘉微(300474，未评级)：GPU 国内排头兵，图形显控与小型雷达高速成长

- 1) 公司自主研发的 JM5400 打破了国外对我国军用 GPU 领域的禁运，创造了军用 GPU 的国产化条件，公司在此基础上不断研发更为先进的 GPU 产品，有望不断满足高端应用需求。

- 2) 公司图形显控模块产品国内技术领先，小型专用化雷达产品逐渐成熟，在核心器件国产化的背景下，有望迎来高速增长。

通富微电(002156，未评级)：收购 AMD 封测资产，切入 GPU 封装和测试领域

- 1) 公司在大基金的支持下，并购 AMD 封测资产并开设合肥、苏通新厂，产能成倍扩大，同时积极布局汽车电子、高性能 CPU 等新兴领域，竞争力显著提升。
- 2) 公司通过收购 AMD 苏州、槟城厂，切入 GPU 等人工智能芯片封装和测试领域，在人工智能芯片封测领域抢占先机。

富瀚微(300613，未评级)：深耕视频监控芯片，智能安防前景广阔。

- 1) 公司深耕视频监控芯片领域，产品布局符合产业发展趋势，受益于安防行业快速发展，公司迅速成长，并与海康等安防龙头企业建立了深度合作关系。
- 2) 公司 14 年底与赛蓝科技合作推出新一代 WIFI 智能可视门铃，积极推进技术升级，基于深度学习算法，将人工智能适当介入安防视频监控，包括对目标分类、区分识别等不同场景。

## 风险提示

人工智能芯片研发不及预期，人工智能芯片技术难度较高，研发周期较长，研发投入较大，相关机构或企业存在研发进度不及预期，以及因无法承受财务压力而退出市场的风险。

下游需求不及预期，人工智能目前已在汽车、机器人、家居和服务器等领域得到应用，下游应用领域存在对人工智能芯片需求不及预期的风险。



## 信息披露

依据《发布证券研究报告暂行规定》以下条款：

发布对具体股票作出明确估值和投资评级的证券研究报告时，公司持有该股票达到相关上市公司已发行股份1%以上的，应当在证券研究报告中向客户披露本公司持有该股票的情况，

就本证券研究报告中涉及符合上述条件的股票，向客户披露本公司持有该股票的情况如下：

截止本报告发布之日，东证资管仍持有海康威视 (002415.SZ) 股票达到相关上市公司已发行股份1%以上。

提请客户在阅读和使用本研究报告时充分考虑以上披露信息。

起点学院，产品、运营技能提升平台

## 分析师声明

每位负责撰写本研究报告全部或部分内容的研究分析师在此作以下声明：

分析师在本报告中对所提及的证券或发行人发表的任何建议和观点均准确地反映了其个人对该证券或发行人的看法和判断；分析师薪酬的任何组成部分无论是在过去、现在及将来，均与其在本研究报告中所表述的具体建议或观点无任何直接或间接的关系。

## 投资评级和相关定义

报告发布日后的 12 个月内的公司的涨跌幅相对同期的上证指数/深证成指的涨跌幅为基准；

### 公司投资评级的量化标准

买入：相对强于市场基准指数收益率 15%以上；

增持：相对强于市场基准指数收益率 5%~15%；

中性：相对于市场基准指数收益率在-5%~+5%之间波动；

减持：相对弱于市场基准指数收益率在-5%以下。

未评级 —— 由于在报告发出之时该股票不在本公司研究覆盖范围内，分析师基于当时对该股票的研究状况，未给予投资评级相关信息。

暂停评级 —— 根据监管制度及本公司相关规定，研究报告发布之时该投资对象可能与本公司存在潜在的利益冲突情形；亦或是研究报告发布当时该股票的价值和价格分析存在重大不确定性，缺乏足够的研究依据支持分析师给出明确投资评级；分析师在上述情况下暂停对该股票给予投资评级等信息，投资者需要注意在此报告发布之前曾给予该股票的投资评级、盈利预测及目标价格等信息不再有效。

### 行业投资评级的量化标准：

看好：相对强于市场基准指数收益率 5%以上；

中性：相对于市场基准指数收益率在-5%~+5%之间波动；

看淡：相对于市场基准指数收益率在-5%以下。

未评级：由于在报告发出之时该行业不在本公司研究覆盖范围内，分析师基于当时对该行业的研究状况，未给予投资评级等相关信息。

暂停评级：由于研究报告发布当时该行业的投资价值分析存在重大不确定性，缺乏足够的研究依据支持分析师给出明确行业投资评级；分析师在上述情况下暂停对该行业给予投资评级信息，投资者需要注意在此报告发布之前曾给予该行业的投资评级信息不再有效。

备注：起点学院 学员收集资料于网络，版权为原作者所有。



起点学院，产品、运营 技能提升平台。