



FPGA硬件加速 为Hadoop+深度学习插上翅膀

邬刚 / Ericwu
ericwu@speed-clouds.com
18603078033



加速云
SPEED CLOUDS



专业引领梦想



我们致力于提供硬件加速的专业技术解决方案、产品和服务。加速云硬件加速产品可以广泛应用于数据中心、云计算、机器视觉、深度学习、仿真、金融、高性能计算等领域。



加速云
SPEED CLOUDS



加速云的技术应用场景





加速云
SPEED CLOUDS



专业的解决方案和产品

解决方案

- ◆ Hadoop平台加速方案
- ◆ 深度学习加速方案
- ◆ 网络安全加速方案
- ◆ 仿真加速和硬件在环仿真方案
- ◆ Spark平台加速方案
- ◆ 金融行业数据加速方案
- ◆ 大数据存储加速方案
- ◆ 高性能计算方案
- ◆ 机器视觉加速方案

01



硬件平台

- ◆ PCIe加速卡
- ◆ 加速模块
- ◆ VPX加速平台
- ◆ 机器视觉加速套件
- ◆ 高密度服务器
- ◆ 定制加速服务器

02



软件及IP

- ◆ 集成开发环境 (SIDE)
- ◆ 深度学习库IP
- ◆ OPENBLAS库IP
- ◆ 压缩解压缩IP
- ◆ 纠删码IP
- ◆ 加解密IP
- ◆ FFT IP
- ◆ 各种定制IP

03





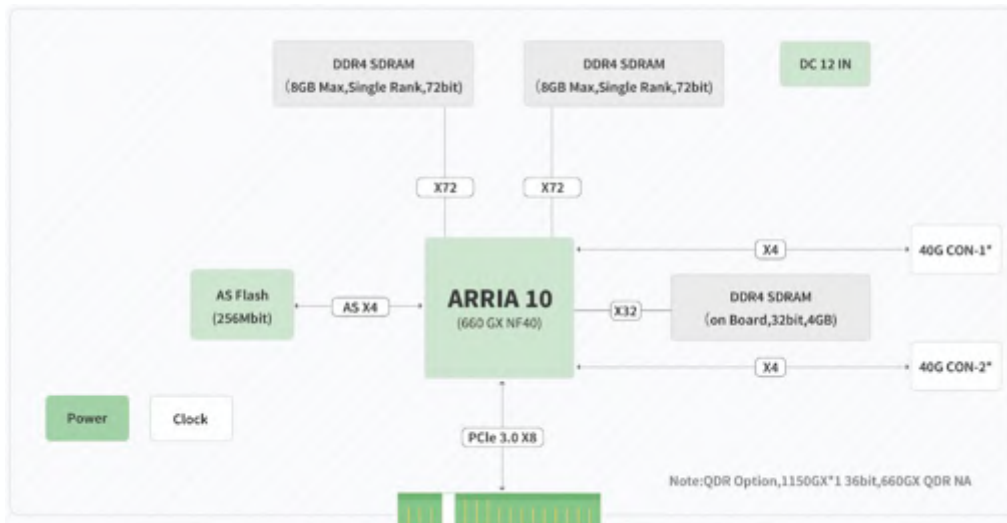
加速云
SPEED CLOUDS



硬件加速平台-PCIe加速板卡 SC-OPM



支持OpenCL 开发, 支持SCIDE开发环境
物理尺寸: 半高半长 (56*167mm)



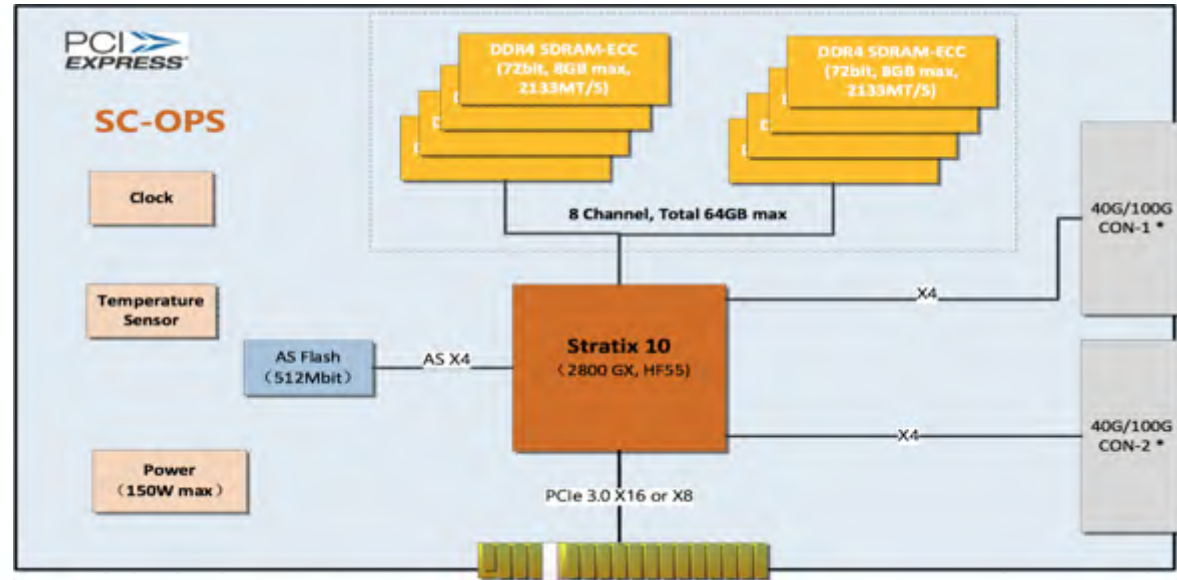
- 采用Altera Arria 10 GX660器件,集成 660k LE和1.5T FLOPS单精度浮点处理能力
- 板卡支持2 个QSFP 电口可以配置为4*1GE/4*10GE/SRIO/Infiniband
- 板卡支持 PCIe 3.0 8Lane 接口,访问带宽为64Gbps
- 板载2组DDR4 2066MHz 72bit颗粒和1组DDR4 2066MHz 32bit颗粒, 支持20GB内存访问带宽共330Gbps
- 支持OpenCL 开发, 支持SCIDE开发环境
- 物理尺寸: 半高半长 (56*167mm)



加速云
SPEED CLOUDS



硬件加速平台-PCIe加速板卡 SC-OPS



- 最新14nm 工艺FPGA S10,逻辑容量2800K, 9.2TFLOPS单精度浮点,18.4TFLOPS16位定点
- 8个内存控制器, 支持2400MHz 72bit DDR4 (ES 2133MHz), 最大支持64GB内存
- PCIe3.0 8lane或16lane (H-Lite支持)
- 支持2个40GE或100GE接口 (H-Lite支持)
- 标准全高3/4长 (112*250mm)
- 正在研发阶段, 预计2017年4月份出样机, 6月份量产



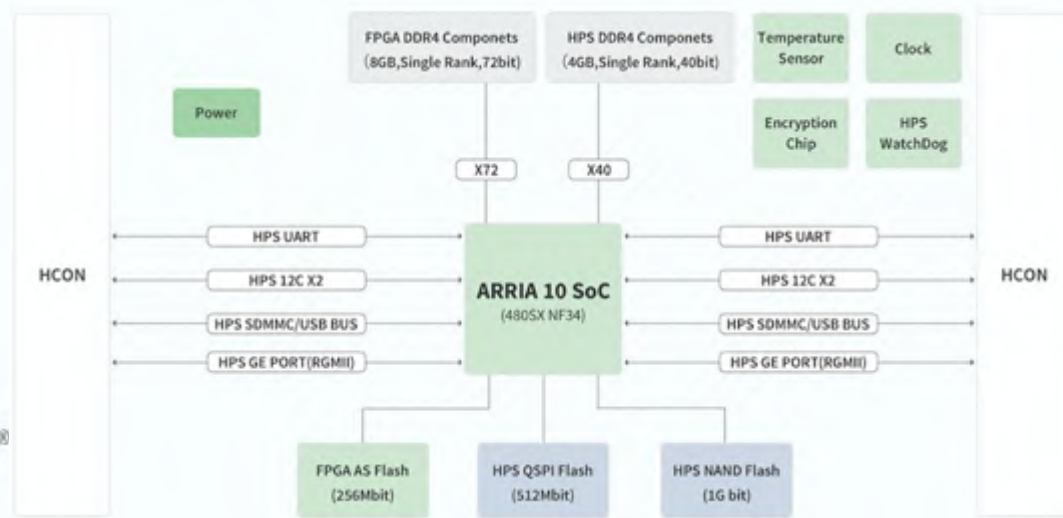
加速云
SPEED CLOUDS



硬件加速平台-核心计算模块 SC-IMB



Dual-core ARM® Cortex®
1.2T/750G FLOPS



- 采用Altera Arria 10 SX480/270器件,集成480k/270k LE和1.2T/750G FLOPS 单精度浮点处理能力
- 集成Dual-core ARM® Cortex®-A9 MPCore™ processor, 最高主频1.5GHz
- FPGA侧支持24个10.3125G SerDes,48对LVDS,101个GPIO管脚,都连接到高速连接器上
- FPGA侧支持1个DDR4 2066MHz 内存控制器,板载最大8GB内存颗粒
- ARM侧支持1个DDR4 2066MHz 内存控制器,板载最大4GB内存颗粒
- ARM侧支持1个千兆网口和USB2.0接口, 2*IIC,1*UART都连接到高速连接器上
- 支持OpenCL 开发,支持SCIDE开发环境
- 物理尺寸: (70*120mm)



支持3D Framing格式输入输出
最高支持4k*2k/60帧分辨率



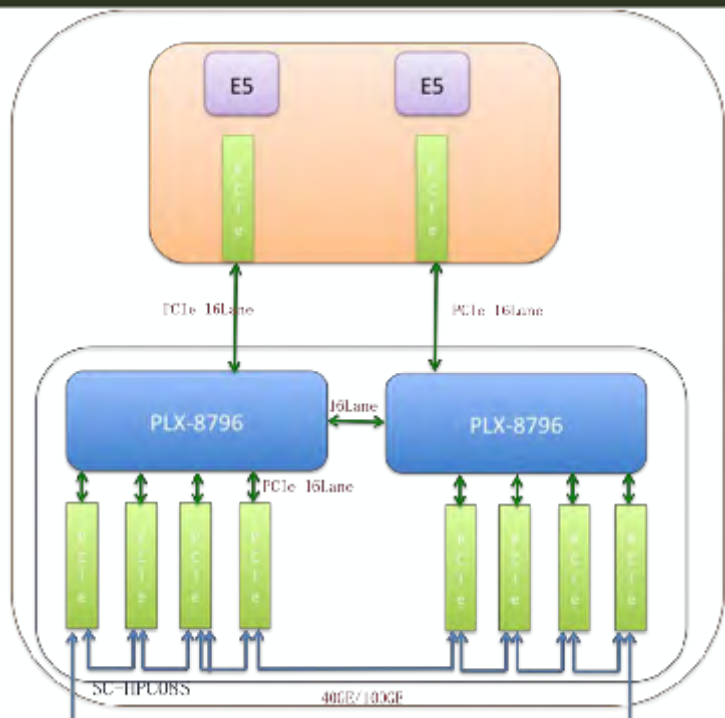
- 支持SC-IMB加速模块作为核心板，提供最高480kLE和1.2TFLOPS 单精度浮点处理能力，集成Dual-core ARM® Cortex®-A9 MPCore™ processor，最高主频1.5GHz
- 支持2路HDMI 2.0输入 + 1路HDMI 2.0输出，最高支持4k*2k/60帧分辨率，支持3D Framing格式输入输出
- 支持1路LVDS 输出，最高支持4k*2k/60帧分辨率
- 支持1路USB3.0接口支持OTG模式
- ARM 侧支持1路USB2.0接口和1路串口，1路千兆以太网电口
- FPGA侧支持1路千兆以太网电口，1路SATA3.0接口
- 支持1路HSMC扩展接口（支持10对高速SerDes+22 IO）



加速云
SPEED CLOUDS



高密度异构计算平台 SC-HPC08S



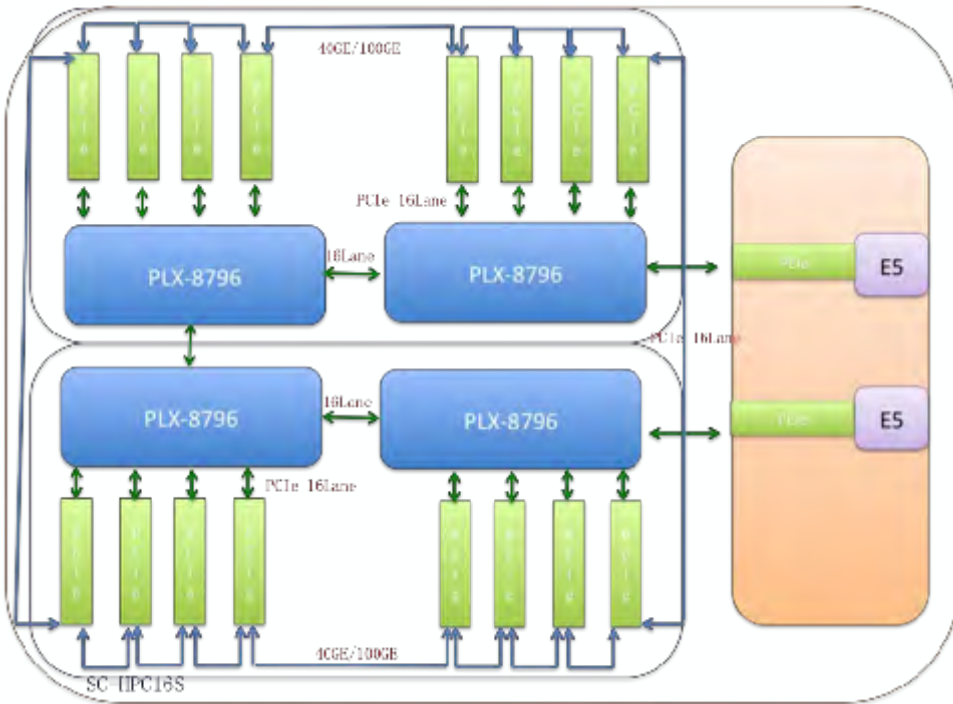
- 19英寸4U标准机箱
- 支持8个PCIe 3.0 16Lane 槽位
- 220V 3+1冗余电源，最大支持4000W
- 两个E5 处理器，每个支持32GB内存
- 2*PCIe 3.0 16Lane (20GB) 互联带宽 (X86和异构计算卡之间)
- 支持FPGA加速卡 (SC-OPM/SC-OPS), 支持GPGPU卡 (K20/K40/K80/M4/M60/P4/P40)
- 在FPGA加速卡时卡间支持40GE/100GE (接口也可以配置为SRIO/infiniband/Serdes/) 互联
- 超高性能功耗比 (在插入SC-OPS时，整个系统支持73.6TFLOPS单精度浮点，功耗为1200W)
- 可以广泛应用于数字信号处理、高性能计算、深度学习等领域



加速云
SPEED CLOUDS



高密度异构计算平台 SC-HPC16S



- 19英寸5U标准机箱
- 支持16个PCIe 3.0 16Lane 槽位
- 220V 3+1冗余电源，最大支持8000W
- 两个E5 处理器，每个支持32GB内存
- 2*PCIe 3.0 16Lane (20GB) 互联带宽 (X86和异构计算卡之间)
- 支持FPGA加速卡 (SC-OPM/SC-OPS), 支持GPGPU卡 (K20/K40/K80/M4/M60/P4/P40)
- 在FPGA加速卡时卡间支持40GE/100GE (接口也可以配置为SRIO/infinband/Serdes/) 互联
- 超高性能功耗比 (在插入SC-OPS时, 整个系统支持150 TFLOPS单精度浮点, 功耗为2000W)
- 可以广泛应用于数字信号处理、高性能计算、深度学习等领域



加速云
SPEED CLOUDS



为什么选择FPGA

	A10 GX660	M4	S10 GX2800	P40
性能 TFLOPS	1.5TFLOPS	2.2TFLOPS	9.2TFLOPS	12TFLOPS
功耗 W	33W	75W	100W	250W
性能/功耗 GFLOPS/W	45GFLOPS/W	29GFLOPS/W	92GFLOPS/W	48GFLOPS/W



加速云
SPEED CLOUDS



解决方案和产品优势

支持更多应用场景

更高的性能功耗比

更高的性能功耗比可以节省
数据中心运营成本



01

02

03

04

对随机操作、位操作和串行算法很好支持可以适应更多应用场景，提高系统性能；低功耗和小型化可以满足更多对功耗小型化有要求的场景

更灵活快速部署

通过加速云FPGA深度学习编译器具有快速迁移相应深度学习算法到加速卡上。利用局部可重构技术可以远程快速部署，满足数据中心云化需求



更高性价比

低功耗可以降低系统运营成本，高集成度可以降低建设成本，更高的性能功耗比使的整体系统获得更高性价比



加速云
SPEED CLOUDS



FPGA加速深度学习

海量数据



01

02



算法

03



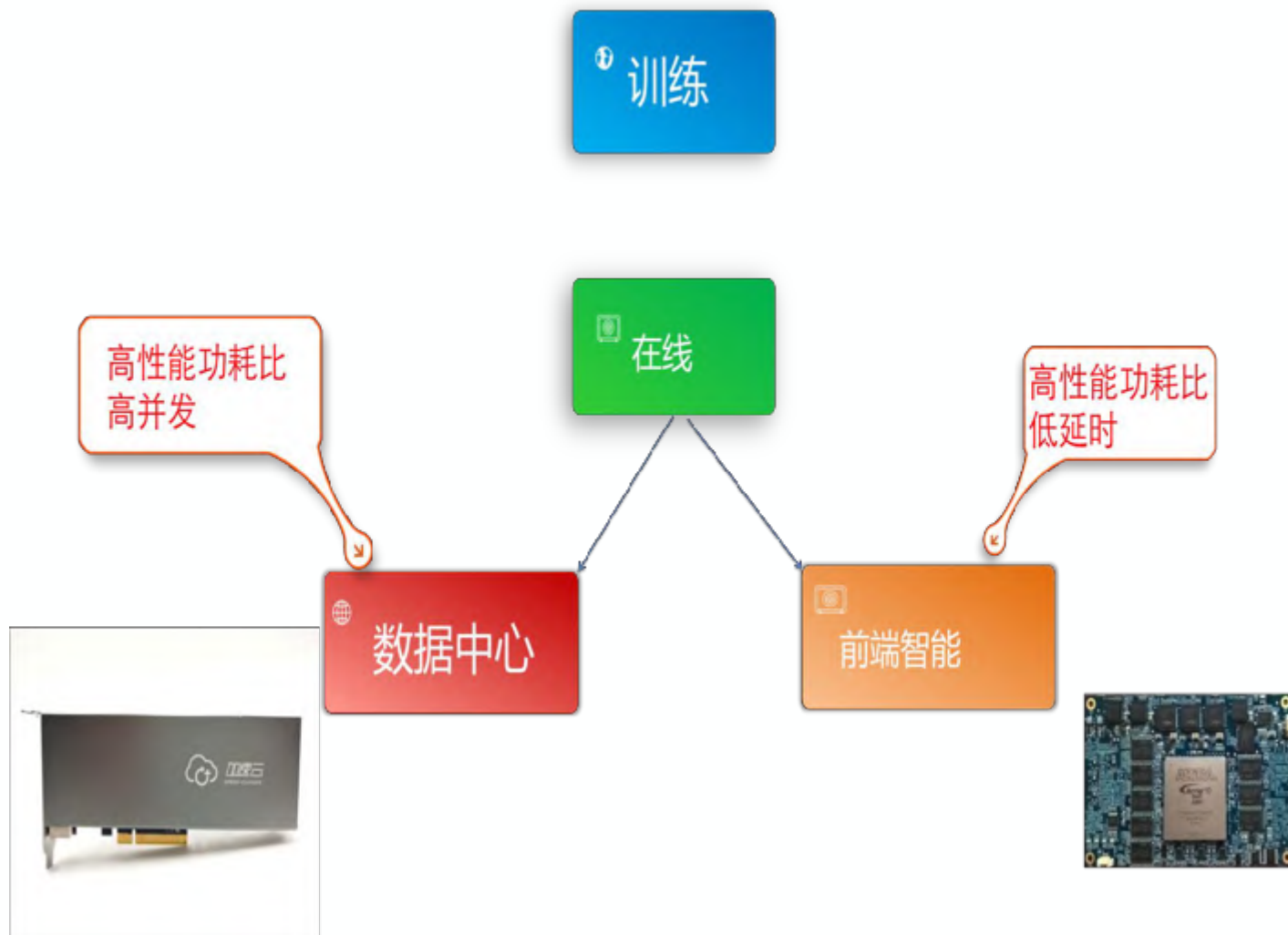
计算能力



加速云
SPEED CLOUDS



深度学习的应用模式





加速云
SPEED CLOUDS



高性能灵活的RTL级加速库



深度学习库
FDNN

参数可配的深度学习基础库：卷积、池化、全连接、非线性函数
参数可配置的CNN/DNN/RNN库，可以兼容CAFFE/TensorFlow模型数据
常见各种模型：VGG16, Goolenet, Lenet, Yolo, SSD, Resnet, Faster-RCNN
参数可配置的深度学习训练库：除CNN/DNN/RNN库外，后向更新算法、随机初始化算法、SGD算法



高性能计算库
FBlas

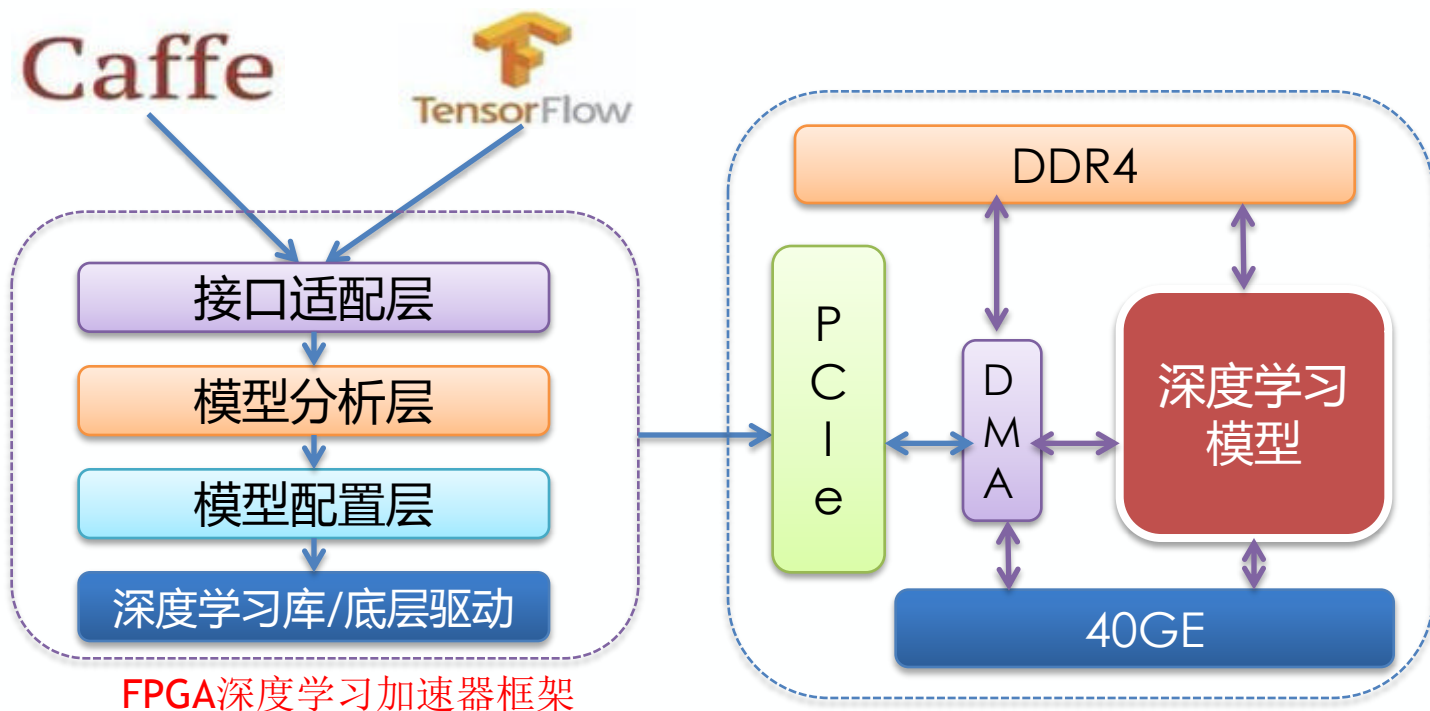
参数可配的OpenBlas库Level2/3:
矩阵乘、矩阵分解、矩阵求逆
线性方程求解、微分方程求解
三角函数、非线性函数、超越函数
傅里叶运算
接口兼容OpenBlas库接口



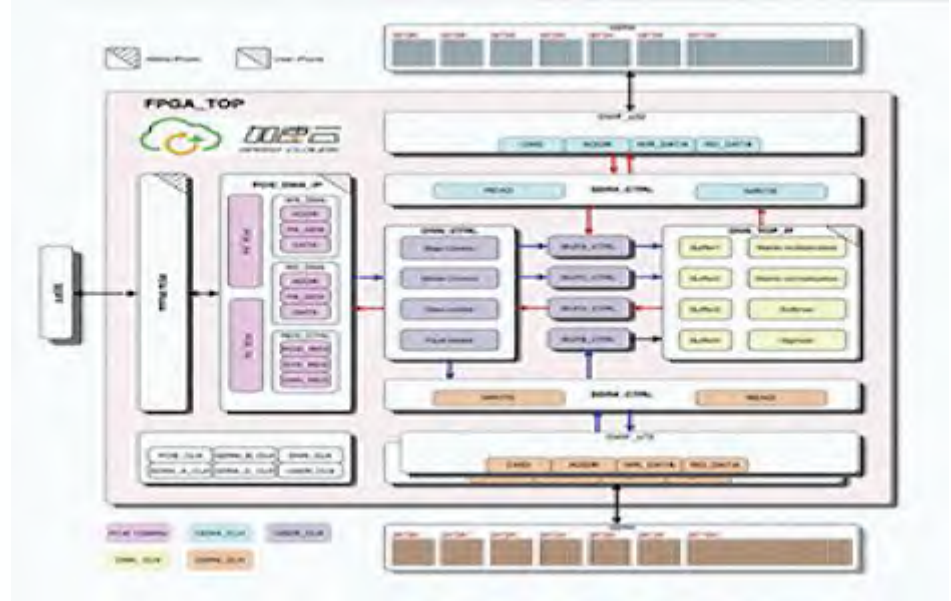
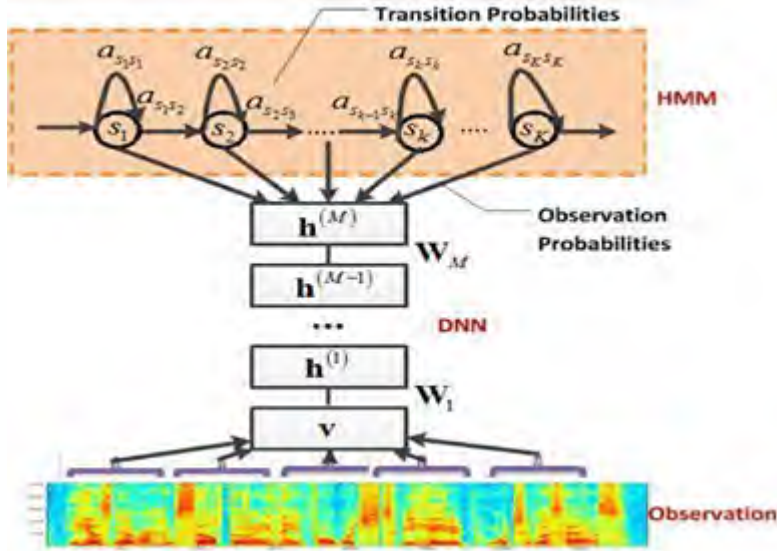
加速云
SPEED CLOUDS



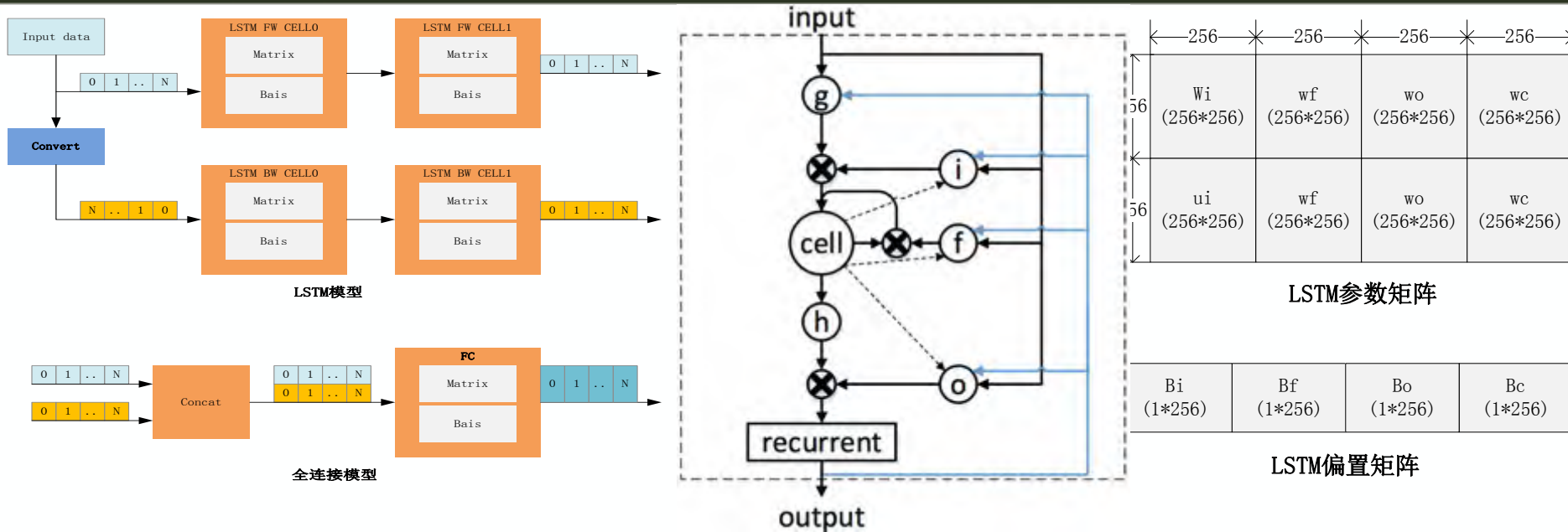
“所设即所加速”的深度学习加速器



通过加速云高性能FPGA深度学习加速器，可以方便和Caffe TensorFlow集成，快速将Caffe TensorFlow训练的模型和数据运行到加速云的FPGA加速卡上（SC-OPM）取得很好的加速比。也可以和加速云高密度异构计算平台配合实现高效的深度学习训练。



- 采用SC-OPM加速卡（半高半长：56*167mm）
- Altera Arria 10 GX660器件,集成 660k LE和1.5T FLOPS 单精度浮点处理能力
- 整体网络为7层，总运算量为84M单精度浮点，激活函数为sigmoid/softmax
- 各层网络参数可以软件配置下载
- 单卡可以实现60路（单精度浮点）语音识别声学模型，8ms全部完成，功耗33瓦
- 单卡可以实现120~150路（16位定点）语音识别声学模型，8ms全部完成
- 采用SC-HPC08S/SC-HPC16S高密度异构计算平台可以实现更高密度语音加速池方案（单系统实现720~1440路语音识别），整体系统最高不超过900W功耗



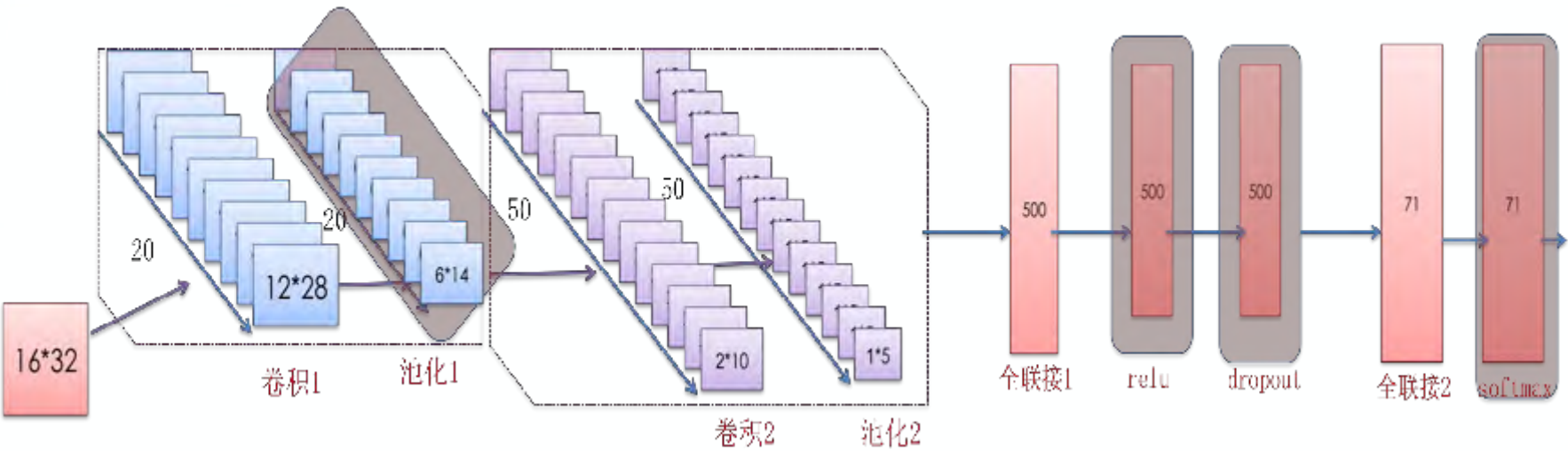
- 采用SC-OPM加速卡（半高半长：56*167mm）
- Altera Arria 10 GX660器件,集成 660k LE和1.5T FLOPS 单精度浮点处理能力
- 四层LSTM+1层全连接，各层网络参数可以软件配置下载
- 可以实现40000T/S的流量，延时超低，数据长度可以混合长度
- 单卡只有33W



加速云
SPEED CLOUDS



深度学习加速解决方案



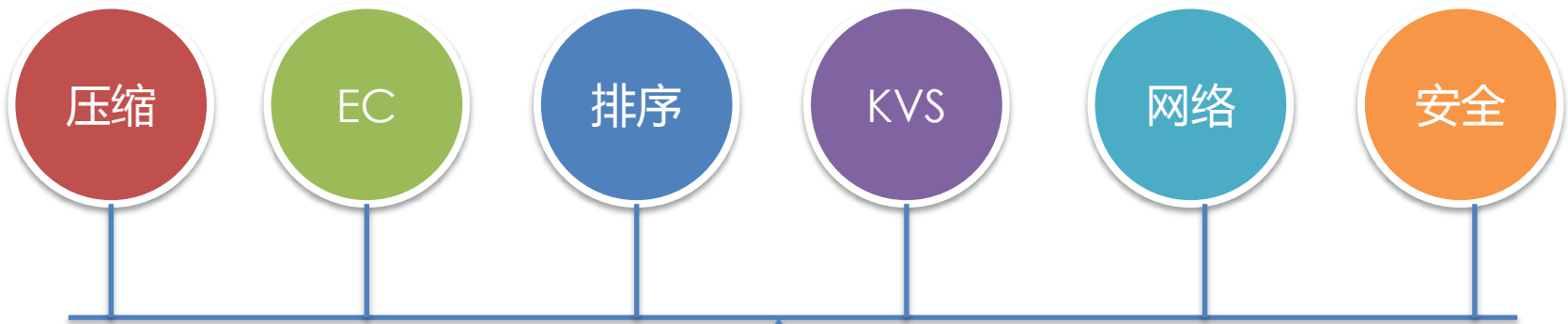
- 采用SC-OPM加速卡（半高半长：56*167mm）
- Altera Arria 10 GX660器件,集成 660k LE和1.5T FLOPS 单精度浮点处理能力
- 2层卷积和2层全连接
- 27400帧/S，单卡只有33W

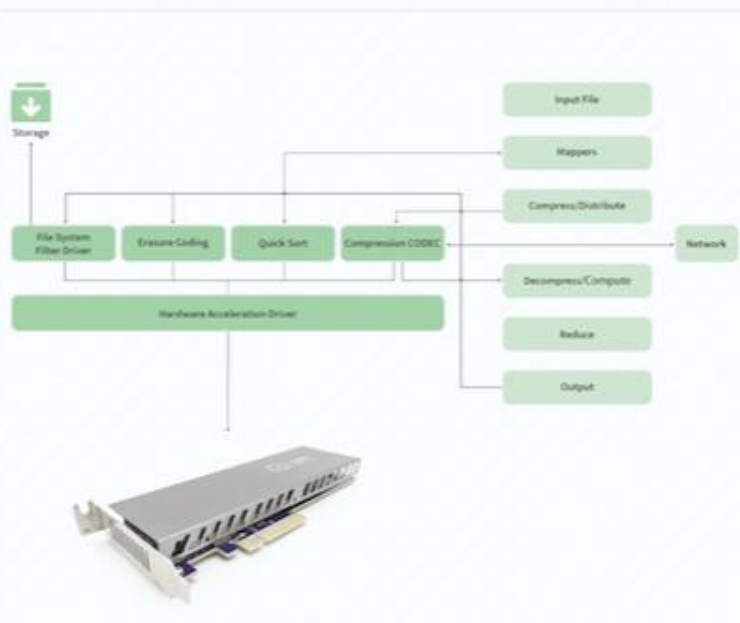
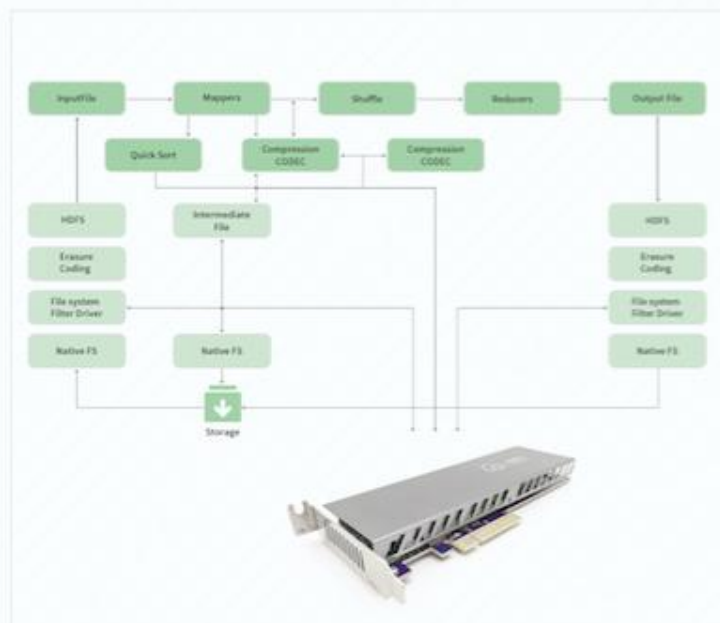


加速云
SPEED CLOUDS



FPGA加速Hadoop





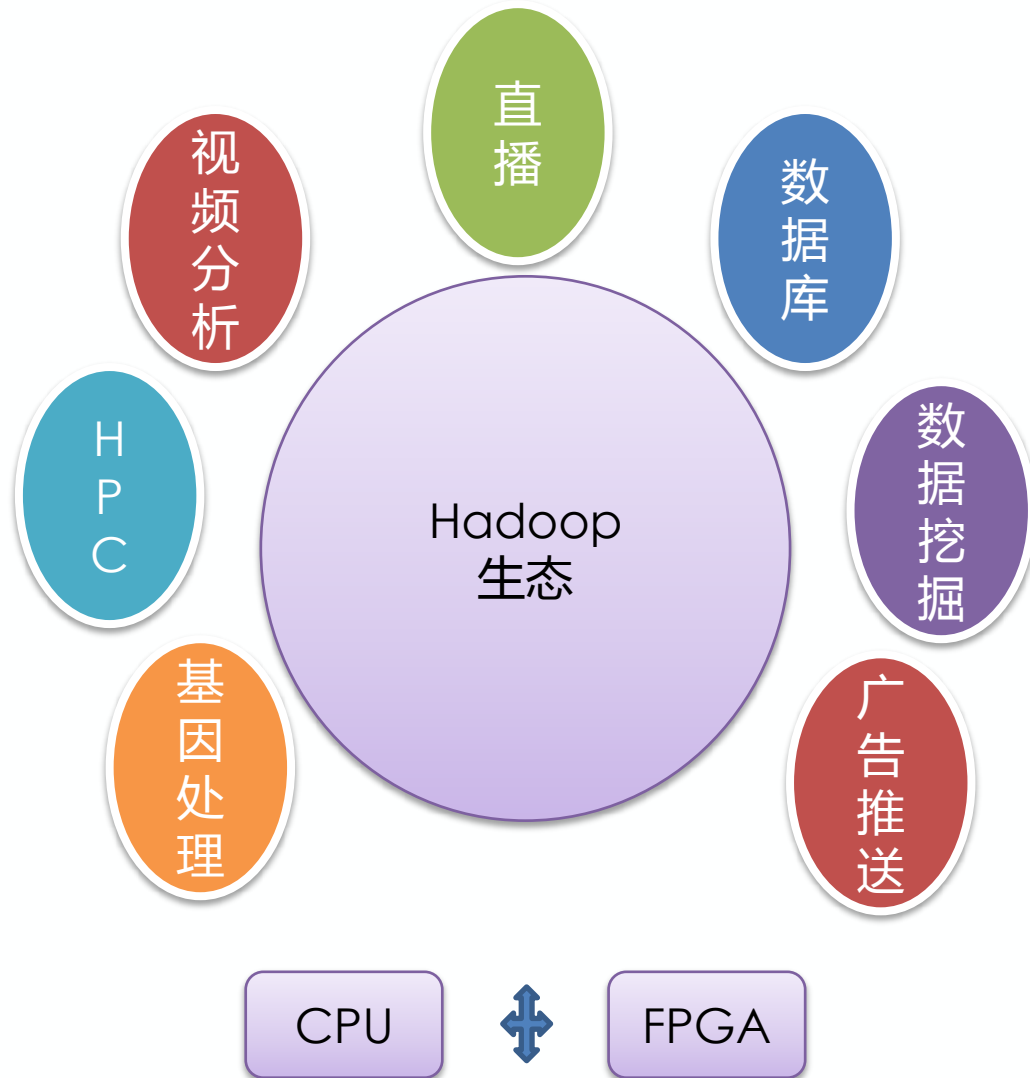
- 实现对Hadoop的相关算法硬件卸载加速，比如压缩解压缩、纠错码（RS）、键值存储（KVS）等
- 可以兼容现有Hadoop大部分版本及LINUX相关版本，兼容EXT3、EXT4、XFS文件系统
- 部署简单，只需要插入加速卡，安装驱动补丁即可
- 降低建设成本和运营成本，可以降低至原有20%左右的成本
- 以压缩解压缩为例：可以是存储节点数降低一半，存储密度提高61%，建设成本降低27%，运营成本降低20%
- 该方案硬件可以编程，支持客户二次开发（支持Opencl和RTL级语言），随时更改算法，比如针对虚拟机可以提供虚拟交换加速
- 采用SC-HPC08S/SC-HPC16S高密度异构计算平台可以实现更高密度压缩解压缩，纠错码加速池方案（单系统实现16~32GB的压缩解压缩性能或24~48GB的纠错码性能），整体系统最高不超过900W功耗



加速云
SPEED CLOUDS



FPGA加速Hadoop应用





加速云
SPEED CLOUDS



基于FPGA硬件加速的Hadoop融合架构

