# INTELIGÊNCIA ARTIFICIAL NA ARQUITETURA INTEL

Igor Freitas

Intel

# Legal Notices and Disclaimers

# Optimization Notice

# Agenda

- Inteligência Artificial (IA) na Arquitetura Intel

- Big Data, IA e Computação de Alta Performance

- Casos de sucesso nacionais

- Tutoriais, treinamentos e informações sobre IA na Arquitetura Intel

# (intel®) AI PORTFOLIO

**SOLUTIONS**

Data Scientists · Technical Services · Reference Solutions

**PLATFORMS** — Intel® AI DevCloud · Intel® Deep Learning System‡ · (intel) Saffron REASONING

**TOOLS** — Intel® Deep Learning Studio‡ · Intel® Deep Learning Deployment Toolkit† · Intel® Computer Vision SDK† · Intel® Movidius™ Software Development Kit (SDK)

**FRAMEWORKS** — TensorFlow* · Caffe* · mxnet* · BigDL on Spark* · neon · Caffe2* · PYTORCH* · CNTK*† · PaddlePaddle*‡

**LIBRARIES** — Intel® MKL/MKL-DNN, clDNN, DAAL, Intel Python Distribution, etc. — DIRECT OPTIMIZATION · Intel® nGraph™ Library† · CPU Transformer† · NNP Transformer‡ · Other

**TECHNOLOGY** — (intel XEON inside) · (intel NERVANA inside) · (intel ARRIA 10 inside) · (intel ATOM inside) · (intel CORE i7 inside) · Iris Graphics · (intel Movidius) — END-TO-END COMPUTE · SYSTEMS & COMPONENTS

†Beta available
‡ Future
*Other names and brands may be claimed as the property of others.

# INTEL AI PLATFORMS

## INTEL® AI DEVCLOUD

Use your existing Intel® Xeon® processor-based cluster
–OR–
Get 4-weeks access to our cluster for FREE including 200GB storage, pre-configured libraries & frameworks

## BETA
## INTEL® DEEP LEARNING SYSTEM*

DL Studio
+
Frameworks
+
Libraries
+
Processors

Enterprise-centric "turnkey" deep learning stack available via rackable on-premise system

## INTEL® SAFFRON COGNITIVE SOLUTION

Collaborative AI decision system for fraud detection, prescriptive maintenance, churn analysis, root cause analysis & more

*Available from select OEMs in 2018

# INTEL AI TOOLS

## FUTURE
### INTEL® DEEP LEARNING STUDIO



Enterprise customer tool to compress full DL development cycle; coming to the Intel® Deep Learning System

## BETA
### INTEL® DEEP LEARNING DEPLOYMENT TOOLKIT



Trained Model → Convert & Optimize → Inference on HW Target

Facilitates optimized inference deployment models trained using the following frameworks: TensorFlow, Caffe, or MXNet

## BETA
### INTEL© COMPUTER VISION SDK



DL deployment kit -PLUS- OpenCV* and OpenVX* support for deep learning-based computer vision on CPU, IPU, GPU & FPGA

### INTEL® MOVIDIUS™ SDK



Intel® Movidius™ Vision Processing Units (VPU) software development kit for inference deployment

*Other names and brands may be claimed as the property of others.*

# INTEL AI FRAMEWORKS

## Popular DL Frameworks are now optimized for CPU!

**CHOOSE YOUR FAVORITE FRAMEWORK**

| TensorFlow™ * | Caffe * | mxnet * | BigDL FOR Apache Spark * | neon™ |
|---|---|---|---|---|

See installation guides at  ai.intel.com/framework-optimizations/

*More under optimization:* ☕ Caffe2 *  PYT❂RCH *  Microsoft CNTK *  PaddlePaddle *  *and others to be enabled via Intel® nGraph™ Library*

# INTEL AI LIBRARIES

## DIRECT OPTIMIZATION

**MKL-DNN**
Open-source optimized deep neural network functions for new frameworks

**clDNN**
Open-source optimized deep neural network functions for Intel GPUs

**DAAL**
Data Analytics Acceleration Library for analytics and machine learning

**Intel Python Distribution**
Optimized distribution of most popular & fastest growing language for machine learning

## INTEL® NGRAPH™ LIBRARY

BETA



Translates participating deep learning framework compute graphs into hardware-optimized executables for many different targets
(CPU, GPU, NNP, FPGA, VPU, etc.)

# INTEL AI COMPUTE

**GENERAL AI**

Mainstream AI

+

Flexible Acceleration

**TRAINING**

**DATA CENTER/WORKSTATION**

Mainstream Training

+

Intensive Training ‡

**INFERENCE**

**DATA CENTER/WORKSTATION**

Mainstream Inference

+

Real-time Inference

**EDGE/GATEWAY**

Mainstream Inference

+

Higher Inference Throughput

Vision 1-20W

Speech/Audio 1-100+mW

Autonomous driving

Custom Inference

**DEEP LEARNING**

‡ Future
All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.

# INTEL® XEON® PROCESSOR PLATFORM PERFORMANCE

## Hardware plus optimized software

### INFERENCE THROUGHPUT

Up to

## 198x

Intel® Xeon® Platinum 8180 Processor
higher Intel optimized Caffe GoogleNet v1 with Intel® MKL
inference throughput compared to
Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

Inference and training throughput uses FP32 instructions

### TRAINING THROUGHPUT

Up to

## 127x

Intel® Xeon® Platinum 8180 Processor
higher Intel Optimized Caffe AlexNet with Intel® MKL
training throughput compared to
Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

Optimized
Frameworks

+

Optimized Intel®
MKL Libraries

## Deliver significant AI performance with hardware and software optimizations on Intel® Xeon® Scalable Family

Up to 191X Intel® Xeon® Platinum 8180 Processor higher Intel optimized Caffe Resnet50 with Intel® MKL inference throughput compared to Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe
Up to 93X Intel® Xeon® Platinum 8180 Processor higher Intel optimized Caffe Resnet50 with Intel® MKL training throughput compared to Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

# FAST & EFFICIENT DL SCALING ON CPU

Intel® - SURFsara* Research Collaboration - Multi-Node Intel® Caffe ResNet-50 Scaling Efficiency on 2S Intel® Xeon® Platinum 8160 Processor Cluster



Chart: Speedup vs Nodes. **HIGHER IS BETTER**. Lines: Ideal (black), SURFsara: MareNostrum4/Barcelona (red). Callout: 90% Scaling Efficiency. Y-axis Speedup (0–70), X-axis Nodes (4, 16, 32, 64, 96, 128, 192, 256).

- MareNostrum4 Barcelona Supercomputing Center
- ImageNet-1K
- 256 nodes
- 90% scaling efficiency
- Top-1/Top-5 > 74%/92%
- Batch size of 32 per node
- Global BS=8192
- Throughput: 15170 Images/sec

**Time-To-Train: 70 minutes (50 Epochs)**

## 90% scaling efficiency with up to 74% Top-1 accuracy on 256 nodes

# LEADING AI RESEARCH

Choose a partner on the cutting-edge of AI breakthroughs



*Neuromorphic Computing Test Chip Codenamed "Loihi"*



*Quantum Computing 49-Qubit Test Chip Codenamed "Tangle-Lake"*

*NEW* AI Technologies @Intel Labs

# BIG DATA, INTELIGÊNCIA ARTIFICIAL E COMPUTAÇÃO DE ALTO DESEMPENHO

# Big Data Analytics

## HPC != Big Data Analytics != Inteligência Artificial ?



**HPC**

| | |
|---|---|
| **FORTRAN / C++ Applications**<br>**MPI**<br>*High Performance* | |
| **SLURM**<br>*Supports large scale startup* | |
| **Lustre***<br>*Remote Storage* | |
| **Compute & Memory Focused**<br>*High Performance Components* | |

Server — Storage SSDs — Switch Fabric — Infrastructure

Modelo de Programação

Resource Manager

Sistema de arquivos

Hardware

**Big Data**

| | |
|---|---|
| **Java, Python, Go, etc.***<br>**Applications**<br>**Hadoop***<br>*Simple to Use* | |
| **YARN***<br>*More resilient of hardware failures* | |
| **HDFS*, SPARK***<br>*Local Storage* | |
| **Storage Focused**<br>*Standard Server Components* | |

Server — Storage HDDs — Switch Ethernet — Infrastructure

*Other brands and names are the property of their respective owners.

(intel) AI

# Big Data Analytics
## *HPC em tempo real*



Data

**Big Data +
Small Compute**
e.g. Search, Streaming,
Data Preconditioning

**FAST DATA**

**Big Data + Big Compute**
e.g. *Real-Time* Local Weather
Modeling, Convolutional Neural Nets

**Small Data + Small
Compute**
e.g. Data analysis

**Small Data +
Big Compute**
e.g. Mechanical Design, Multi-physics

Compute

Solution Urgency

### Varied Resource Needs

■ Processor  ■ Memory  ■ Interconnect  ■ Storage

System cost balance

Typical
Big Data
Workloads

Video Survey   Traffic Monitor   Personal Digital Health

System cost balance

Typical HPC
Workloads

High Frequency Trading   Numeric Weather Simulation   Oil & Gas Seismic

# Trends in HPC + Big Data Analytics



**Performance**
- Code Modernization (Vector instructions)
- Many-core
- FPGA, ASICs
- Integrated solutions: Storage + Network + Processing + Memory

**Usability**
- Frameworks
- Libraries + Program. Lang. (E.g. Python)

**Standards**
- Commom Environments
- Portability
- Open

**Business viability**
- Better products
- Easy to mantain HW & SW
- Faster time-to-market
- Lower costs (HPC at Cloud ? )
- Public investments

# High Performance Deep Learning for **FREE** on CPU Infrastructure[1]



| DataFrame | | | | | |
|---|---|---|---|---|---|
| | | | ML Pipelines | | |
| SQL | SparkR | Streaming | MLlib | GraphX | **BigDL** |
| Spark Core | | | | | |

BigDL is a distributed deep learning library for Apache Spark* that can run directly on top of existing Spark or Apache Hadoop* clusters with direct access to stored data and tool/workflow consistency!

*No need to deploy costly accelerators, duplicate data, or suffer through scaling headaches!*

**Feature Parity** with Caffe* and Torch*

**Lower TCO and improved ease of use** with existing infrastructure

Deep Learning on Big Data Platform, Enabling **Efficient Scale-Out**

## software.intel.com/bigdl

[1]Open-source software is available for download at no cost; 'free' is also contingent upon running on existing idle CPU infrastructure where the operating cost is treated as a 'sunk' cost

# Unified Big Data Analytics Platform



## Hadoop & Spark Platform

Machine Leaning

Graph Analytics

SQL

Notebook

Spreadsheet

Batch | Streaming | Interactive | R* | Java* | Python*

Data Processing & Analysis

DataFrame

ML Pipelines

SQL* | SparkR* | Streaming* | Mllib* | GraphX*

Spark Core

Flink* | Storm*

MR* | Giraph*

Resource Mgmt & Co-ordination

YARN* | ZooKeeper*

Data Input

Flume* | Kafka*

Storage

HDFS* | Parquet* | Avro* | Hbase*

How to Run Deep Leaning Workloads Directly on Big Data Platform?

- Integrated with Big Data ecosystem
- Massively distributed, shared-nothing
- Scale-out
- Send compute to data
- Fault tolerance
- Elasticity
- Incremental scaling
- Dynamic resource sharing
- ...

# Centros de Excelência Intel em Inteligência Artificial c/ Startups

*Apoio de P&D no desenvolvimento de soluções inovadoras em IA*

Centros de Excelência em Inteligência Artificial em parceria com a Intel

**AI2BIZ**  **AXONDATA**

- Transferência de conhecimento em IA e HPC (High Performance Computing)
- Apoio no desenvolvimento de protótipos utilizando software livre
- Workshop em IA e HPC
- Acesso a servidores Intel de alto desempenho



**Entendimento do problema**
- Definição do problema
- Apresentação de casos de sucesso
- Como a Intel pode ajudar

**Workshop**
- Treinamento técnico em:
  - IA na prática
  - Melhoria de performance
- Hands-on

**Prova de Conceito / MVP**
- Baseado em código livre
- 30 à 90 dias
- Repositório de código e docs
- Protótipo otimizado p/ produção

**Fase Piloto**
- Suporte p/ testes
- Otimização de performance

**Fase da solução**
- Suporte p/ deploy
- Otimização de performance
- Casos de sucesso
- Press release
- Eventos

# Centros de Excelência em Inteligência Artificial – Intel
## *Casos de sucesso*



**SSERPRO**
Serviço Federal de Processamento de Dados

"Validador Cognitivo de Infrações de Trânsito"

Thiago Oliveira, superintendente de Engenharia de Infraestrutura do SERPRO

✓ Performance 22.5x mais rápida em "Xeon Scalable Processors"

*"…um processamento de multas que antes levava 45 horas agora poderá ser realizado em menos de 2 horas."*

✓ Desenvolvimento do modelo matemático

"Com isso, tivemos uma acurácia de 90% no sistema, além da automação de todo o projeto",
disse Gustavo Rocha, chefe de divisão do SERPRO,"

Intel Python Distribution + Caffe / TensorFlow  otimizados + MKL  +  Técnicas de HPC

libnumactl    kmp_affinity

https://software.intel.com/en-us/articles/boosting-deep-learning-training-inference-performance-on-xeon-and-xeon-phi

# Centros de Excelência em Inteligência Artificial - Intel
## *Provas de Conceito em andamento*

*Automação da Análise de Processos Jurídicos*

- Classificação e agrupamento automático de processos
- Identificação de processos válidos, seguindo critérios do cliente (falta de carimbo, data errada, nome das partes, etc.)
- OCR de documentos mal digitalizados

*Otimização de Deep Learning em Ambiente de Produção*
*Área Financeira – Banco Público*

- Melhoria de performance e acurácia
- Uso da BigDL em cluster Hadoop + Spark
- Do Protótipo ao Produto

*Predição de falhas em Ponto de Vendas (POS)*
*Área Financeira - Setor de Logística e Pagamentos*

- Diminuir custo de manutenção
- Análise histórica das falhas
- Entendimento da eficiência das máquinas
- Eventos externos e internos
- Predição de quando ocorrerão novas falhas
- Análise de 100k terminais

*Predição de falhas em Caixas Eletrônicos (ATM)*
*Área Financeira – Setor de TI*

- Análise histórica das falhas
- Predição de quando ocorrerão novas falhas

# Um pouco do histórico da Intel em HPC no Brasil
## *Casos de Sucesso*

**Oil & Gas - Reservoir Simulator at PETROBRAS**

- Up to 10.5x performance gains in their Reservoir Simulator software[1]

**LNCC - National Laboratory for Scientific Computing**
**Largest HPC cluster in Latin America**

- Up to 30x performance gain in Oil & Gas applications[2]
- Up to 3.4x speedup via AVX (vector instructions)
- Link white-paper

**Health & Life Sciences**

- Up to 11x speedup in Molecular Dynamics – NCC/UNESP & LNCC – white-paper link
  - Xeon only:
    - Original code vs Modernized code: up to 11x speedup
  - Xeon + 1 Xeon Phi (same optimized code)
    - 1.14x speedup
- Article link

**INPE/CPTEC**
**Code Modernization of BRAMS**

- Initial results – white-paper link

Authors:
[1]CENPES team and Gilvan Vieira - gilvandsv@gmail.com
[2]LNCC - Frederico Cabral - fredluiscabral@gmail.com
[3]NCC/UNESP - Silvio Stanzani  silvio.stanzani@gmail.com

26

# TUTORIAIS, TREINAMENTOS E INFORMAÇÕES SOBRE IA NA ARQUITETURA INTEL

# INTEL® AI ACADEMY

For developers, students, instructors and startups

## LEARN



- Online tutorials
- Webinars
- Student kits
- Support forums

## DEVELOP



- Intel Optimized Frameworks
- Exclusive access to Intel® AI DevCloud

## TEACH



- Comprehensive courseware
- Hands-on labs
- Cloud compute
- Technical Support

## SHARE



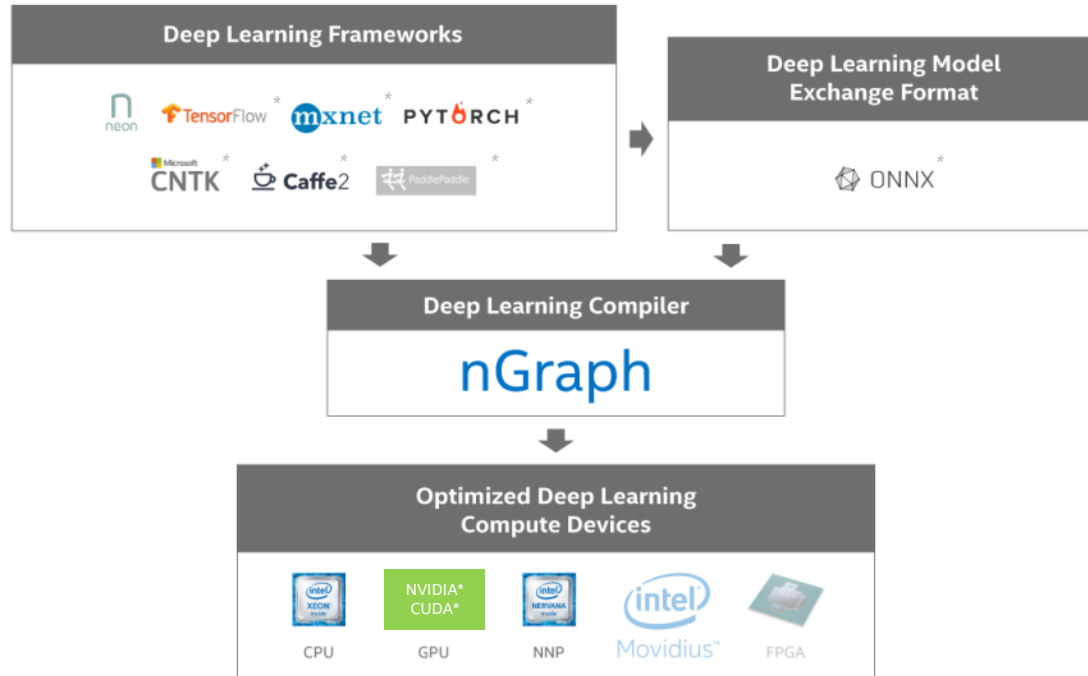- Project showcase opportunities at
- Intel Developer Mesh
- Industry & Academic events

## software.intel.com/ai

# INTEL® NGRAPH™ COMPILER

## Optimized Compute Devices for Neural Networks



*Other names and brands may be claimed as the property of others.*
*All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.*

# INTEL DISTRIBUTION FOR PYTHON

## Advancing Python Performance Closer to Native Speeds

**For developers using the most popular and fastest growing programming language for AI**

### Easy, Out-of-the-box Access to High Performance Python

- Prebuilt, optimized for numerical computing, data analytics, HPC
- Drop in replacement for your existing Python (no code changes required)

### Drive Performance with Multiple Optimization Techniques

- Accelerated NumPy/SciPy/Scikit-Learn with Intel® MKL
- Data analytics with pyDAAL, enhanced thread scheduling with TBB, Jupyter* Notebook interface, Numba, Cython
- Scale easily with optimized MPI4Py and Jupyter notebooks

### Faster Access to Latest Optimizations for Intel Architecture

- Distribution and individual optimized packages available through conda and Anaconda Cloud
- Optimizations upstreamed back to main Python trunk

### software.intel.com/intel-distribution-for-python

# INTEL® MKL-DNN

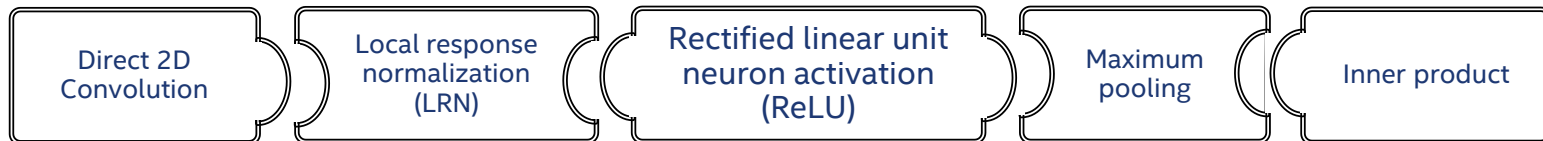## Math Kernel Library for Deep Neural Networks

**For developers of deep learning frameworks featuring optimized performance on Intel hardware**

### Distribution Details

- Open Source
- Apache 2.0 License
- Common DNN APIs across all Intel hardware.
- Rapid release cycles, iterated with the DL community, to best support industry framework integration.
- Highly vectorized & threaded for maximal performance, based on the popular Intel® MKL library.
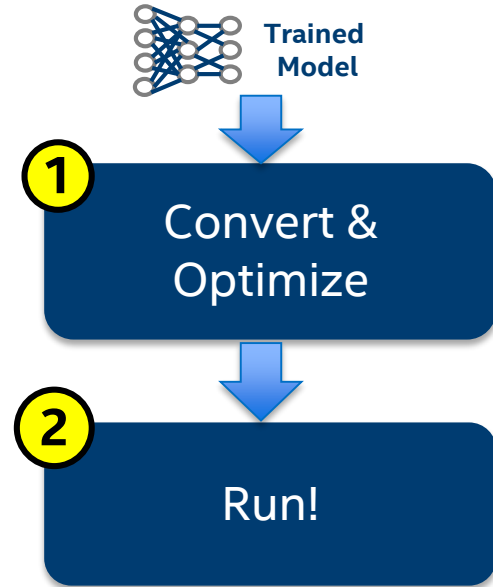
github.com/01org/mkl-dnn

**Examples:** Direct 2D Convolution | Local response normalization (LRN) | Rectified linear unit neuron activation (ReLU) | Maximum pooling | Inner product

# INTEL® DEEP LEARNING DEPLOYMENT TOOLKIT

**BETA Now Available!**

## For developers looking to run deep learning models on the edge

**Trained Model**

**①**

Imports trained models from popular DL framework regardless of training HW

Enhances model for improved execution, storage & transmission

**①** Convert & Optimize

**②**

Optimizes Inference execution for target hardware (computational graph analysis, scheduling, model compression, quantization)

Enables seamless integration with application logic

Delivers embedded friendly Inference solution

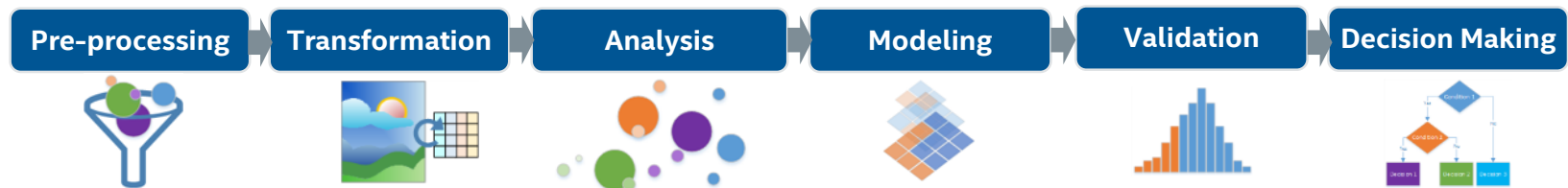**②** Run!

## Ease of use + Embedded friendly + Extra performance boost

# INTEL® DATA ANALYTICS ACCELERATION LIBRARY (INTEL® DAAL)

## High Performance ML and Data Analytics library

**Building blocks for all data analytics stages, including data preparation, data mining & machine learning**

| Pre-processing | Transformation | Analysis | Modeling | Validation | Decision Making |
|:---:|:---:|:---:|:---:|:---:|:---:|

Open Source • Apache 2.0 License

Common Python, Java and C++ APIs across all Intel hardware

Optimized for large data sets including streaming and distributed processing

Flexible interfaces to leading big data platforms including Spark and range of data formats (CSV, SQL, etc.)

# FIND OUT MORE

**LEARN**

**More information at ai.intel.com**

**EXPLORE**

**Use Intel's performance-optimized libraries & frameworks**

**ENGAGE**

**Contact your Intel representative for help and POC opportunities**