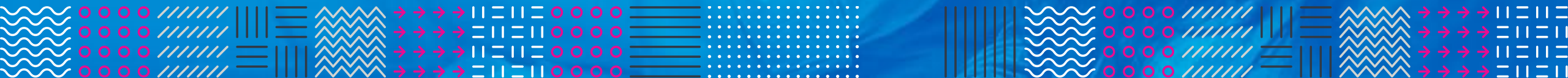




Accelerated ML, instantly

Chris Kachris
CEO, Co-founder



Need for Computing power



“computing power needed to carry out machine learning neural networks is **doubling every 3.5 months.**”

Cliff Young, Google

New ML models need more powerful platform

In 2018, OpenAI found that the amount of computational power used to train the largest AI models had doubled every 3.4 months since 2012.

<https://www.technologyreview.com/s/614700/the-computing-power-needed-to-train-ai-is-now-rising-seven-times-faster-than-ever-before/>

Open AI

www.inaccel.com™

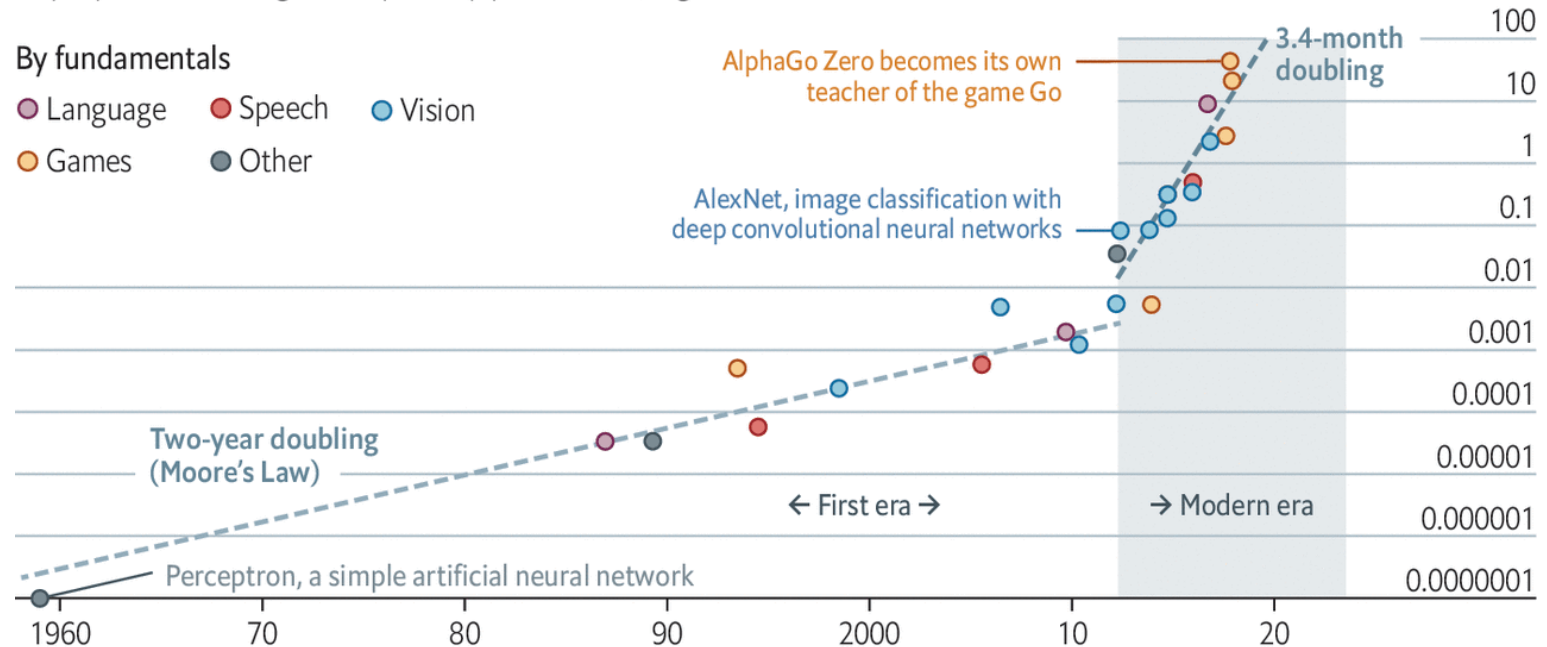
Deep and steep

Computing power used in training AI systems

Days spent calculating at one petaflop per second*, log scale

By fundamentals

- Language
- Speech
- Vision
- Games
- Other



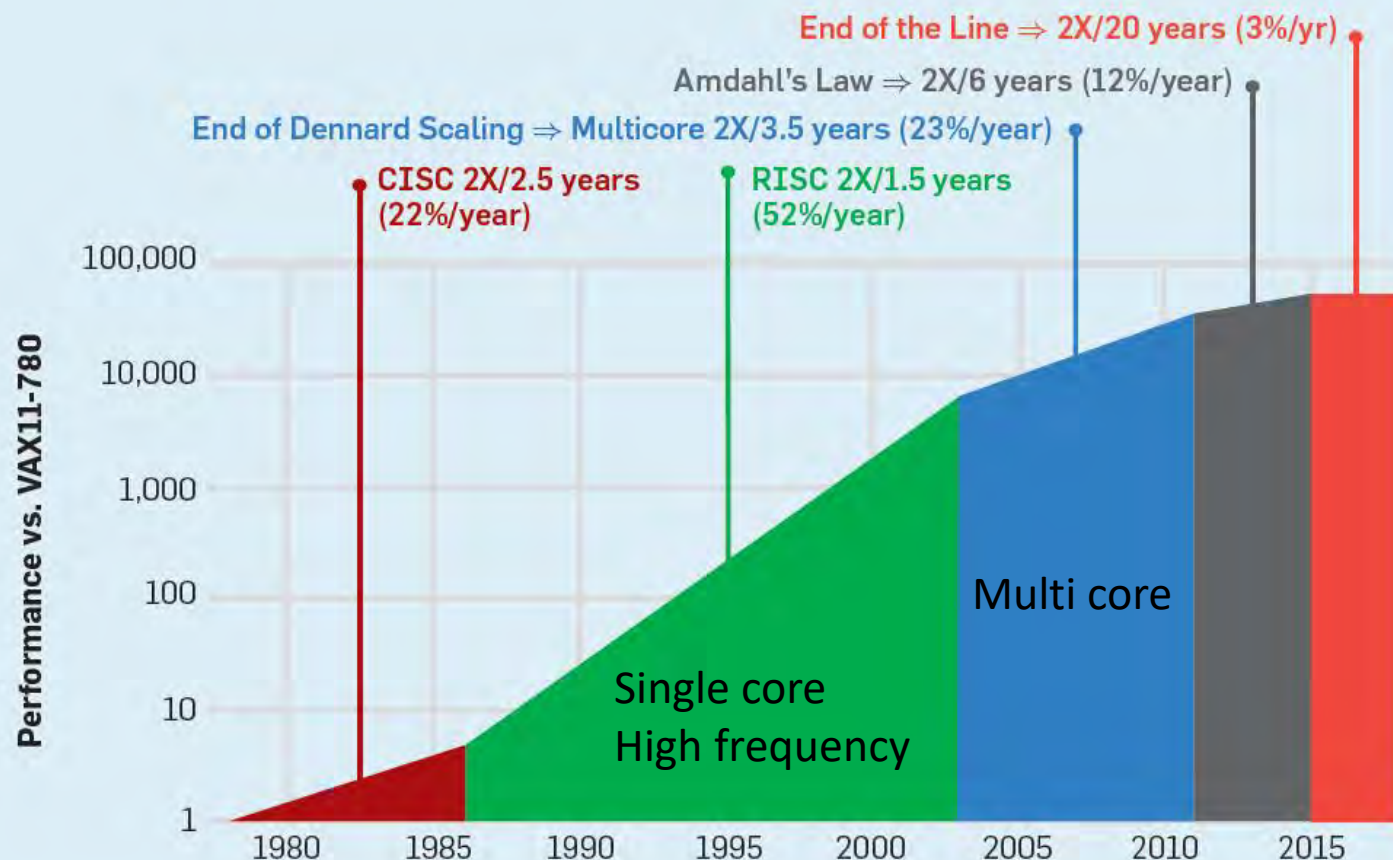
Source: OpenAI

The Economist

*1 petaflop=10¹⁵ calculations

<https://www.economist.com/technology-quarterly/2020/06/11/the-cost-of-training-machines-is-becoming-a-problem>

What's left for faster computing?



CPUs cannot keep increasing the performance as in the past

What's Left?

Only path left is **Domain Specific Accelerators (DSA)**

- Just do a few tasks, but extremely well
- Standard computer for legacy software + accelerators to improve performance per Watt of critical computation

David Patterson, 2019



Next **WAVE** of **COMPUTING**:
Accelerators

Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

A domain-specific architecture for deep neural networks

Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson

Accelerators

Easy to use

CPU



GPU



inaccel

FPGA



Performance



Accelerated Machine Learning



Accelerated Machine Learning

Speedup your applications online using the power of accelerators



<https://inaccel.com/accelerated-data-science/>

What we do: Accelerated Machine Learning



Help companies **speedup** their ML applications by using **accelerators** (FPGAs) seamlessly (**ML as a Service**):



10x – 20x Faster

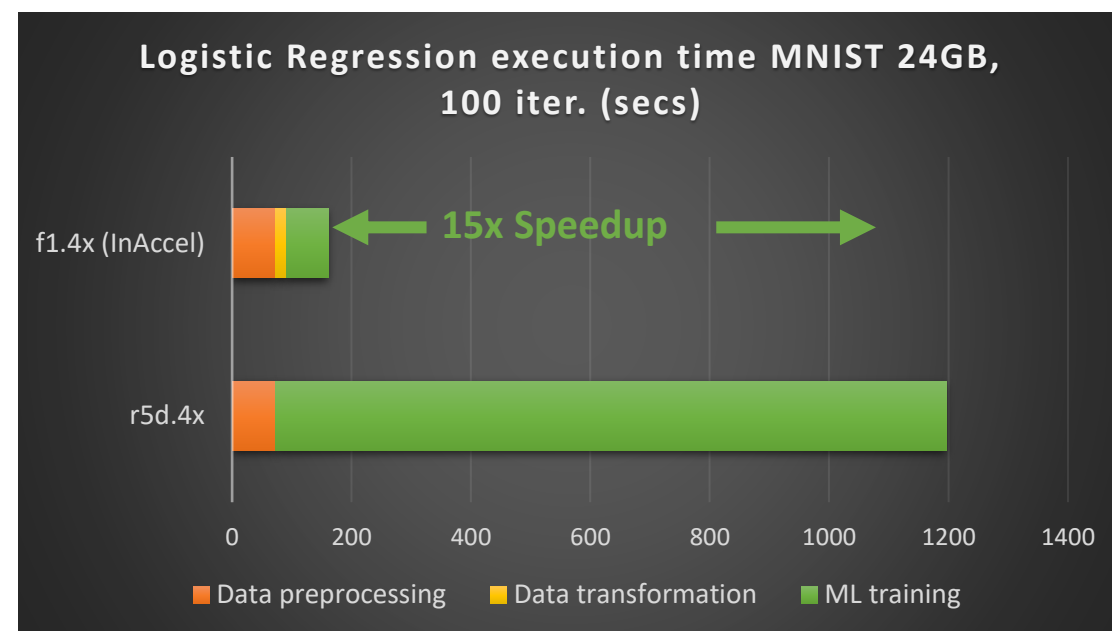


2x Lower cost



Zero code changes

Integrated Accelerated Machine Learning/DNN Platform



<https://www.youtube.com/watch?v=hDLAYNsF39s>

Same tools (Python), Faster results



InAccel Studio is a unique portal for Accelerated Machine learning

A screenshot of a Jupyter Notebook interface. The title bar shows 'KerasExample.ipynb'. The notebook content is titled 'ImageNet Classification with ResNet50'. It includes a text block explaining that the notebook shows how to classify images using InAccel's Keras-like framework. Below this, there are code cells. Cell [1] imports numpy, time, and InAccel's ResNet50 and ImageDataGenerator. Cell [2] loads a ResNet50 model. Cell [3] sets up an ImageDataGenerator and loads images from a directory. Cell [4] performs accelerated inference using the model and prints the duration and FPS. The output of cell [4] shows a duration of 12.287 seconds and an FPS of 1848.498 for 22615 images.

```
[1]: import numpy as np
import time

from inaccel.keras.applications.resnet50 import decode_predictions, ResNet50
from inaccel.keras.preprocessing.image import ImageDataGenerator, load_img
last executed at 2020-06-22 11:06:12 in 11m

[2]: model = ResNet50(weights='imagenet')
last executed at 2020-06-22 11:06:12 in 2m

Accelerated Inference

Then, we load thousands of images specifying the number of batches for every process.

[3]: data = ImageDataGenerator(dtype='float32')
images = data.flow_from_directory('imagenet/', target_size=(224, 224), class_mode=None, batch_size=64)
last executed at 2020-06-22 11:06:18 in 1.3s
Found 22615 Images belonging to 498 classes.

Now, it's time to feed the model with the images and predict their class.

We also measure the performance as the number of Images processed Per Second.

[4]: begin = time.monotonic()
preds = model.predict(images, workers=10)
end = time.monotonic()

print('Duration for: ', len(preds), 'Images: %.3f sec' % (end - begin))
print('FPS: %.2f' % (len(preds) / (end - begin)))
last executed at 2020-06-22 11:06:18 in 11.25s
Duration for 22615 Images: 12.287 sec
FPS: 1848.498
```

- **Familiar tools**



- **Faster results**

- 10x-20x faster
- Zero code changes

<https://inaccel.com/accelerated-machine-learning/>

Instant Acceleration



Just add **inaccel** and enjoy **15x faster** execution

File Edit View Run Kernel Tabs Settings Help | InAccel Cloud light

/ shared / ml /

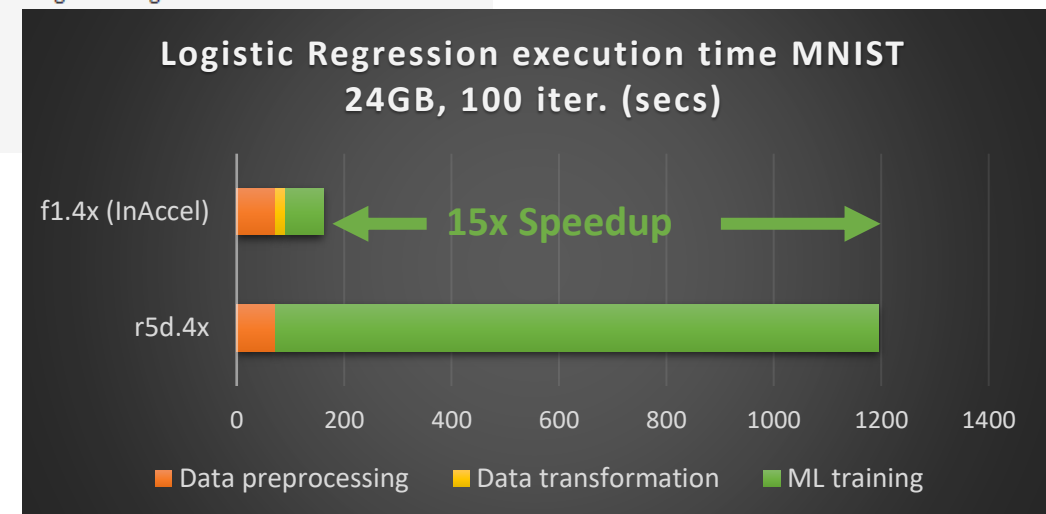
Name	Last Modified
data	12 days ago
img	19 days ago
LogisticRegression.ipynb	a month ago
MLlibPipelines.ipynb	12 days ago
NaiveBayes.ipynb	a month ago
XGBoost.ipynb	19 days ago

Logistic Regression Hyperparameter Tuning using FPGAs

This notebook shows how to train and apply many accelerated sklearn models, with a k-fold cross validation and hyperparameter tuning step.

```
[ ]: from inaccel sklearn.linear_model import LogisticRegression
from sklearn.datasets import fetch_openml
from sklearn.linear_model import LogisticRegression as LogisticRegressionCPU
from sklearn.metrics import accuracy_score
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
```

15x faster ML training



MLPerf



ID	Submitter	System	Processor	#	Accelerator	#	Software	Model Used	Results				
									Task	Image classification			
									Data	ResNet			
									Accuracy (%FP32 ref)	99.00			
									Scenario	Single Stream	Multiple Stream	Server	Offline
									Units	latency in ms	streams	query/s	samples/s
CATEGORY: Available													
0.7-158	Deci	n2-highcpu-16	Intel(R) Xeon(R) Cascade Lake CPU	1			openvino_2020.4.287	resnet		2.13			1,093
0.7-159	Deci	n2-highcpu-2	Intel(R) Xeon(R) Cascade Lake CPU	1			openvino_2020.4.287	resnet		6.77			148
0.7-160	dividiti	AWS g4dn.4xlarge	Intel(R) Xeon(R) Platinum 8259CL CPU @ 2.50GHz	1			OpenVINO 2020	resnet50		2.67			
0.7-161	dividiti	AWS g4dn.4xlarge	Intel(R) Xeon(R) Platinum 8259CL CPU @ 2.50GHz	1			OpenVINO 2020	ssd-resnet34					
0.7-162	InAccel	alveo2.ethz.ch	Intel(R) Xeon(R) Gold 6234	2	Xilinx Alveo U250	1	InAccel Keras 2.3.1.2	resnet		7.00	58		3,026
0.7-163	InAccel	alveo2.ethz.ch	Intel(R) Xeon(R) Gold 6234	2	Xilinx Alveo U250	2	InAccel Keras 2.3.1.2	resnet					6,008

MLPerf inference of ResNet50 on an Alveo U250 cluster



Instant Acceleration by more than 10x



Now we're ready to build a pipeline and fit it. This puts the data through all of the feature processing, model tuning & training we described in a single call.

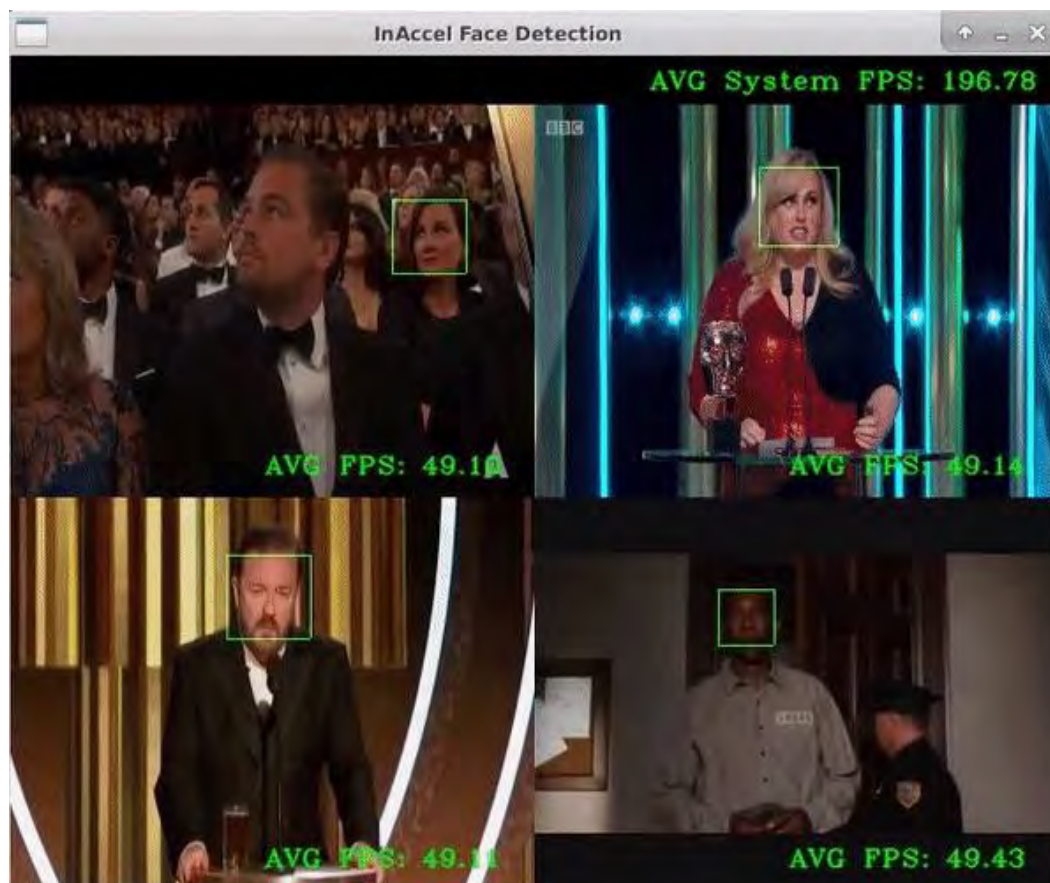
```
In [*]: from pyspark.ml import Pipeline  
  
pipeline = Pipeline(stages=[labelIndexer, featuresScaler, cv, indexToLabel])  
  
%time model = pipeline.fit(train)
```

▼ Apache Spark: 1 EXECUTORS 8 CORES Jobs: 1 RUNNING 6 COMPLETED				
	27 March 16:45			
	38	39	40	41
Jobs:	42			
	3 4 5			

The next cell converts the test set from LibSVM to **Parquet** memory format, in order to serve as our streaming source.

15x
faster

Video analytics: Face detection



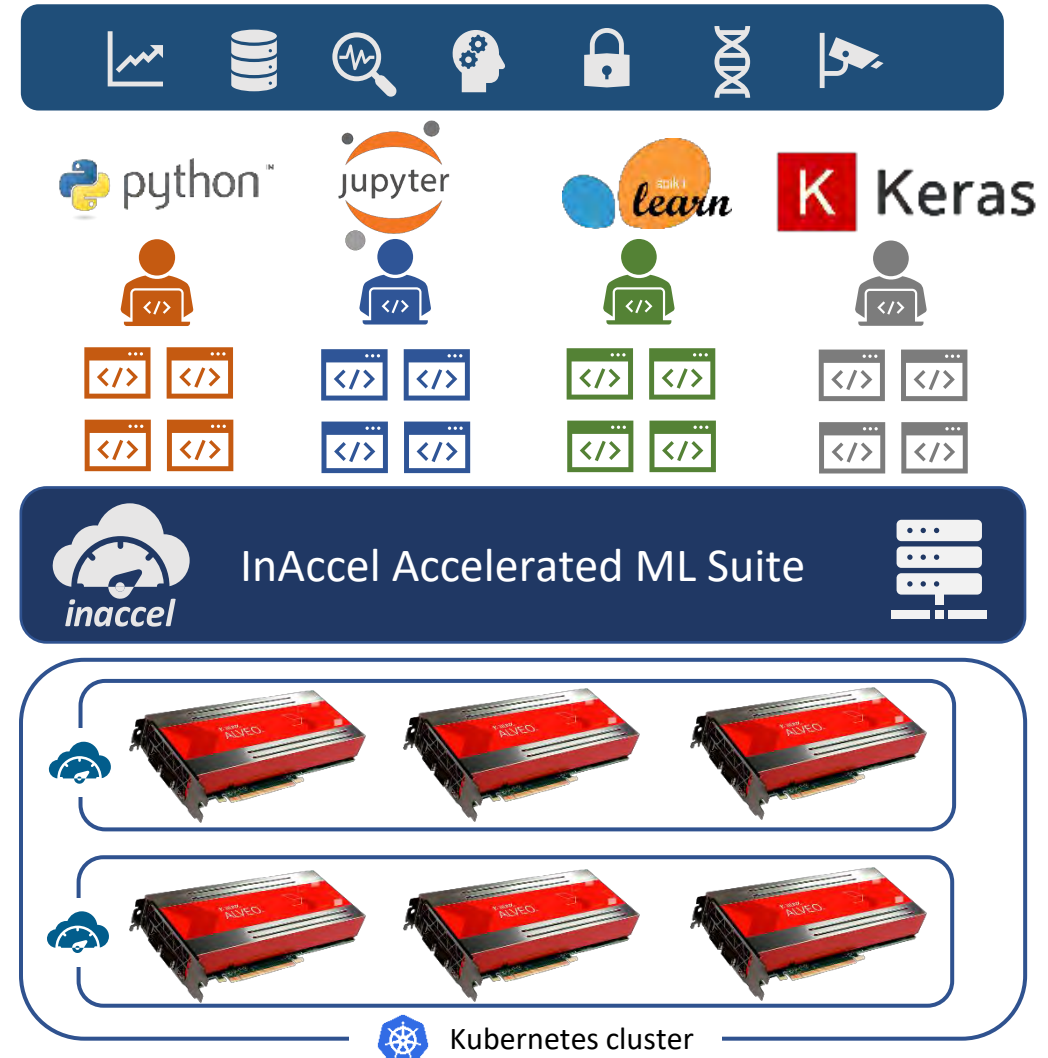
1700fps per server

How we do it



- We leverage the **power of programmable accelerators (FPGA)**
- A unique platform for easy deployment, scaling and resource management of FPGA
- **Pricing model:** Pay-as-you-go (Subscription)

[*https://www.xilinx.com/support/documentation/white_papers/wp504-accel-dnns.pdf](https://www.xilinx.com/support/documentation/white_papers/wp504-accel-dnns.pdf)



Use cases

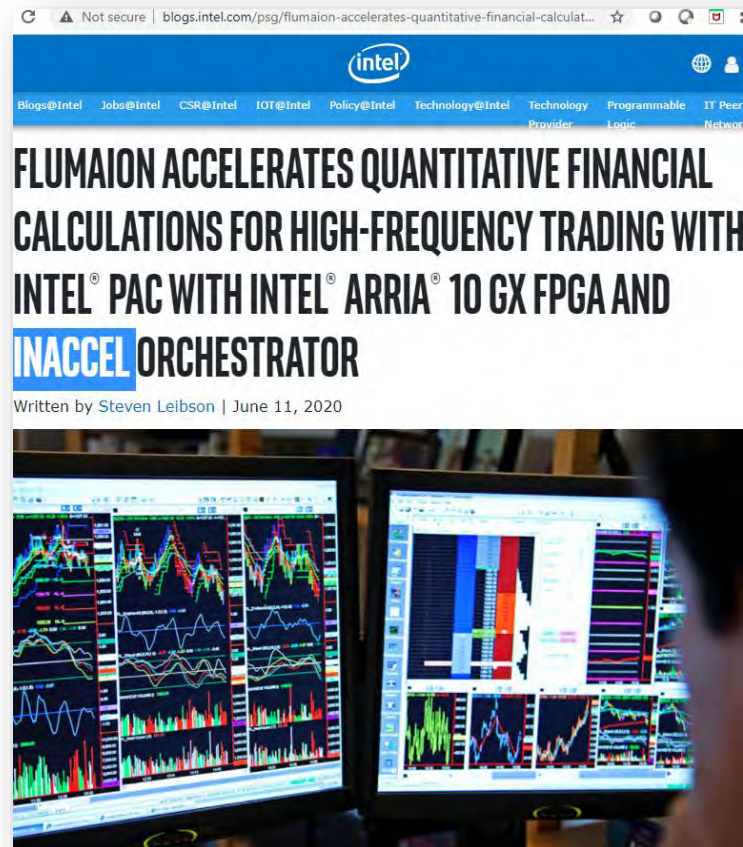


Machine Learning



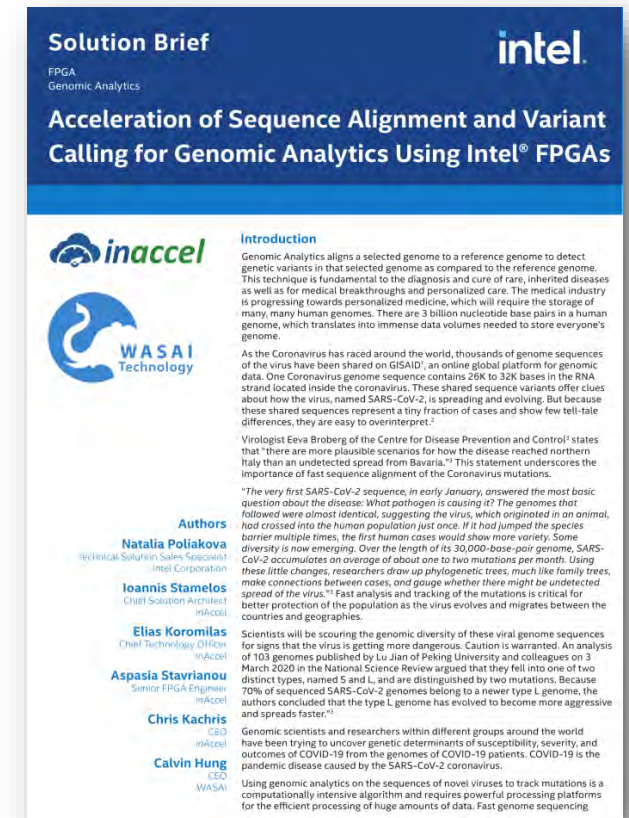
<https://blogs.intel.com/psg/inaccels-accelerated-ml-suite-boosts-spark-ml-performance-by-as-much-as-7x-on-fpga-based-alibaba-cloud-f1-instances/>

Quantitative Finance















<https://blogs.intel.com/psg/flumaion-accelerates-quantitative-financial-calculations/>

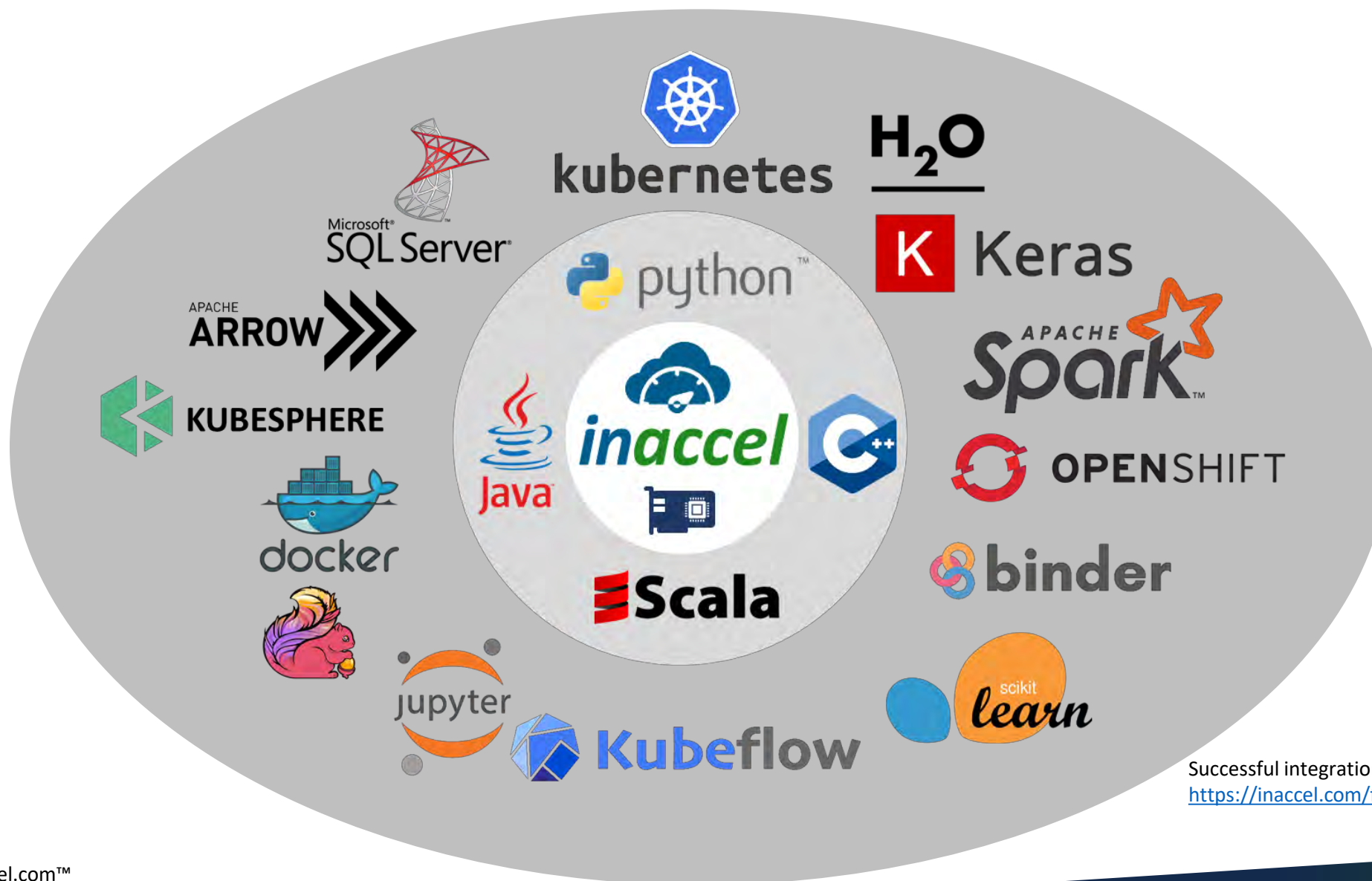
Genomics



On line free Data Science platforms

	GPU	FPGA
		 10x
		
		
		
		
		

Seamless Integration with any framework



Successful integrations with Multiple Frameworks
<https://inaccel.com/fpga-integration-easier-than-ever/>

Market size – FPGAs in Data Centers



InAccel is targeting the Market size for **FPGA Data Center** and **HPC** Accelerators for **application Acceleration**

- Machine learning
- Genomics
- Quantitative Financial
- Analytics
- Databases
- Security
- Vision



2023

TAM: \$500 Million

(Compute acceleration not including smartNIC and storage)

SAM: \$200 Million

(Compute acceleration not including smartNIC and storage)

Users



Flumaion



University of Sarajevo





- (13) **United States**
- (12) **Patent Application Publication**
- Kulkarni et al.**
- (14) Pub. No.: **US 2009/0276599 A1**
- (43) Pub. Date: **Nov. 5, 2009**
- Publication Classification**
- (51) **Int. Cl.**
G06F 12/00 (2006.01)
- (52) **U.S. Cl.**
711/179; 711/112.0
- (57) **ABSTRACT**
- A configurable transactional memory synchronizes transactions from clients. The configurable transactional memory includes a memory buffer and a transactional buffer. The memory buffer includes allocation control and storage; the allocation control is configurable to selectively allocate the storage between a transactional buffer and a data buffer for the data words. The transactional buffer stores state indicating each coalescence of a data word and a client for which the data word is referenced by a write access in the transaction in progress from the client. The transactional buffer organizes data coalescences in the transaction in progress from each client. The completion status is an indicator for collision or aborted for a collision. A collision is occurred that references a data word of the transaction in progress following a write access that references the data word of another transaction in progress from another client.



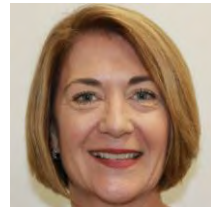
InAccel Team



Advisory board



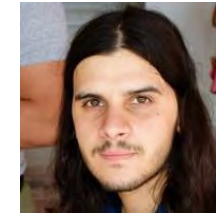
Christos
Makiyama
Founder and
President
at Silicon Planet
Corporation



Genelle Heim
Managing Director at
Grayson Hayden
Group
(Ex-Vice President of
Marketing at Xilinx)



Chris Kachris
CEO, co-founder



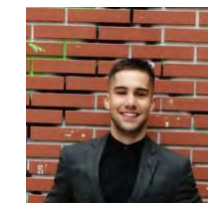
Elias Koromilas
CTO, co-founder



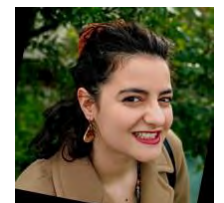
Ioannis Stamelos
COO, co-founder



Vasilis
Amourgianos
FPGA Engineer



Vangelis
Gkiastas
ML Engineer



Aspasia
Stavrianou
DevOps engineer

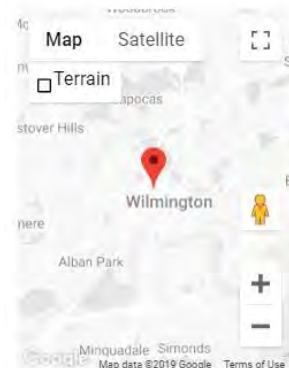
InAccel, Inc. Corporate overview



- Founded in January 2018 (Seed fund: \$600 USD in June'18)
- Registered in Delaware, USA



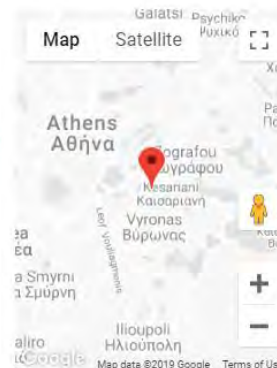
- Membership:



Headquarters

500 Delaware Ave STE 1, #1960
Wilmington, DE 19801
USA

(+1) 408 260 5724



Design Center

Formionos 47
Kesariani 116 33
Athens, Greece

(+30) 211 1825 436





QUESTIONS?



MLconf Online™

THANK YOU!

@inaccel
/in/kachris