



Методы работы с большими данными

Киреев Василий Сергеевич,
к.т.н., доцент

Москва, 2020

Киреев Василий Сергеевич,
к.т.н., доцент

Москва, 2020

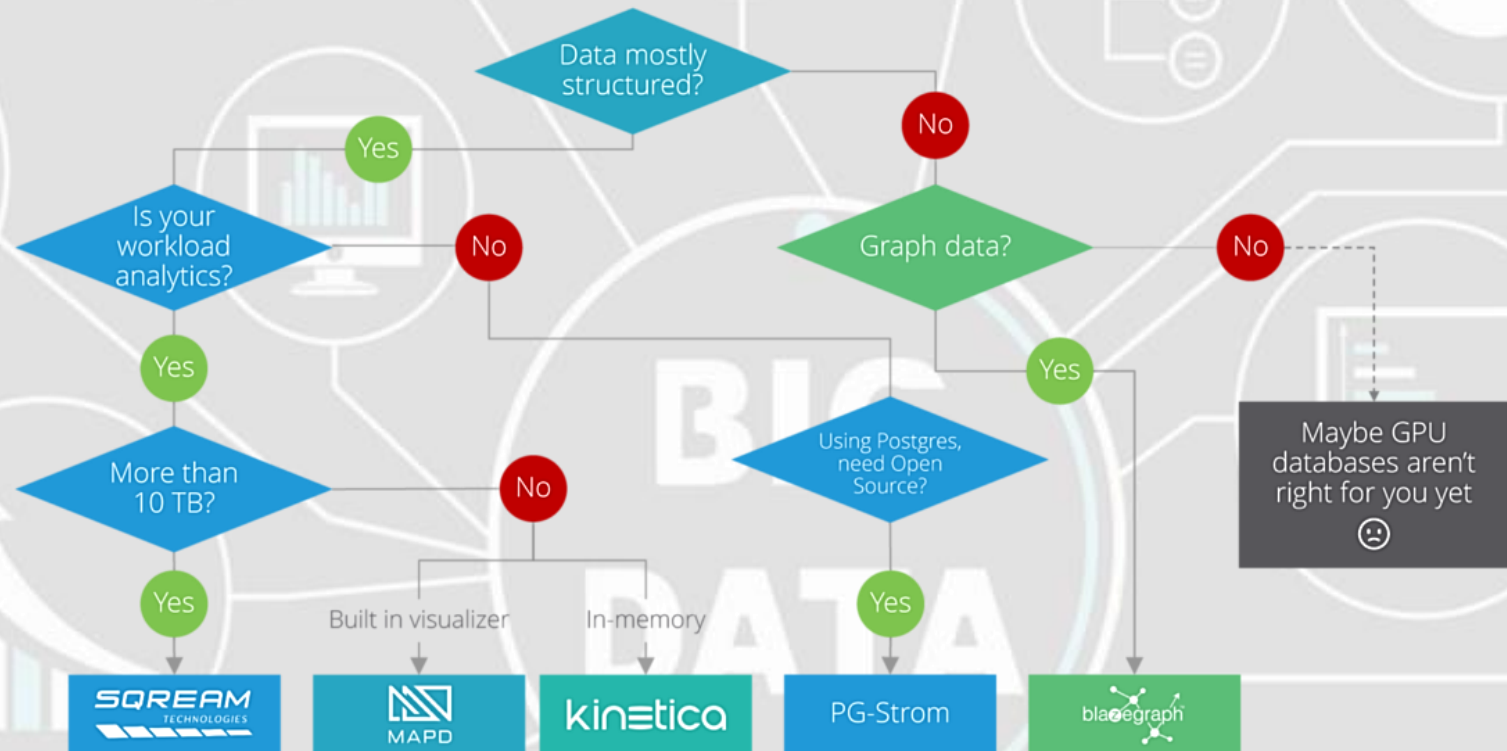
GPU и большие данные

Graphics Processing Unit — это графический процессор широко используемый в настольных и серверных системах. Отличительной особенностью этого устройства является ориентированность на массовые параллельные вычисления. В отличие от графических процессоров архитектура другого вычислительного модуля CPU (Central Processor Unit) предназначена для последовательной обработки данных.

Если количество ядер в обычном CPU измеряется десятками, то в GPU их счет идет на тысячи, что накладывает ограничения на типы выполняемых команд, однако обеспечивает высокую вычислительную производительность в задачах, включающих параллелизм.

GPU в хранении больших данных

Which GPU Database Should I Use?



GPU в обработке больших данных

Hadoop Processing, Reading from disk



Spark In-Memory Processing



Traditional GPU Processing

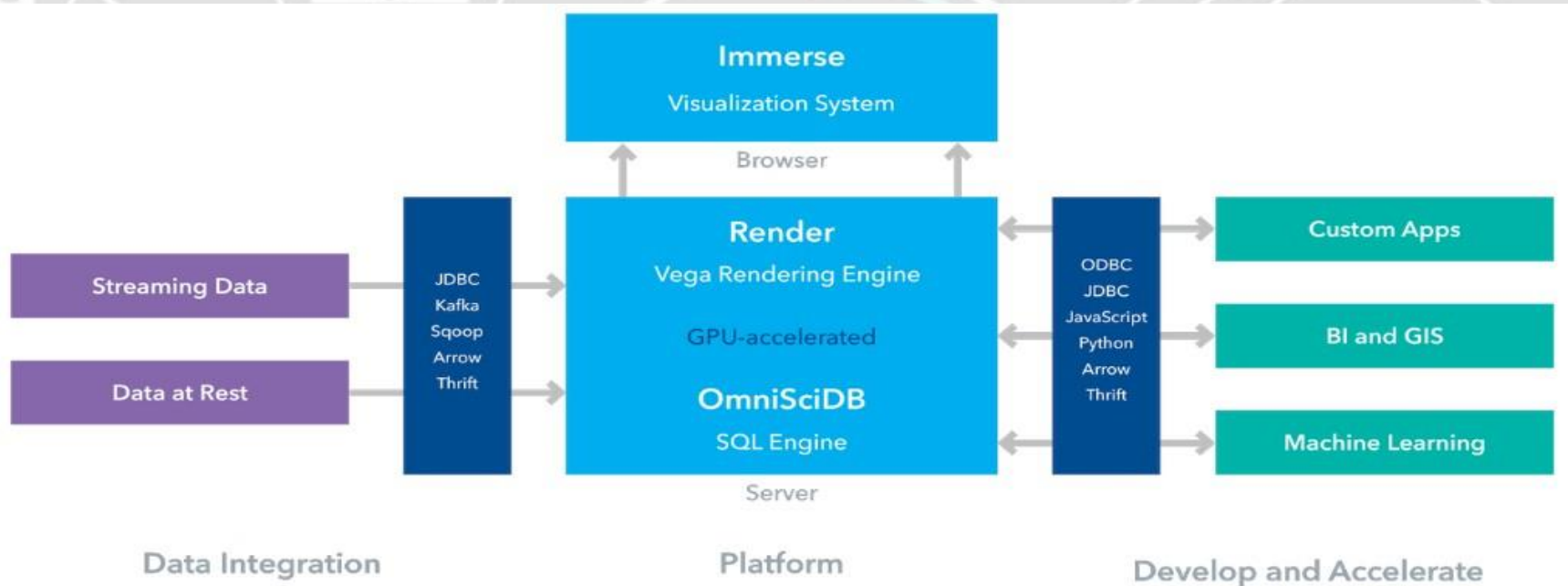


RAPIDS

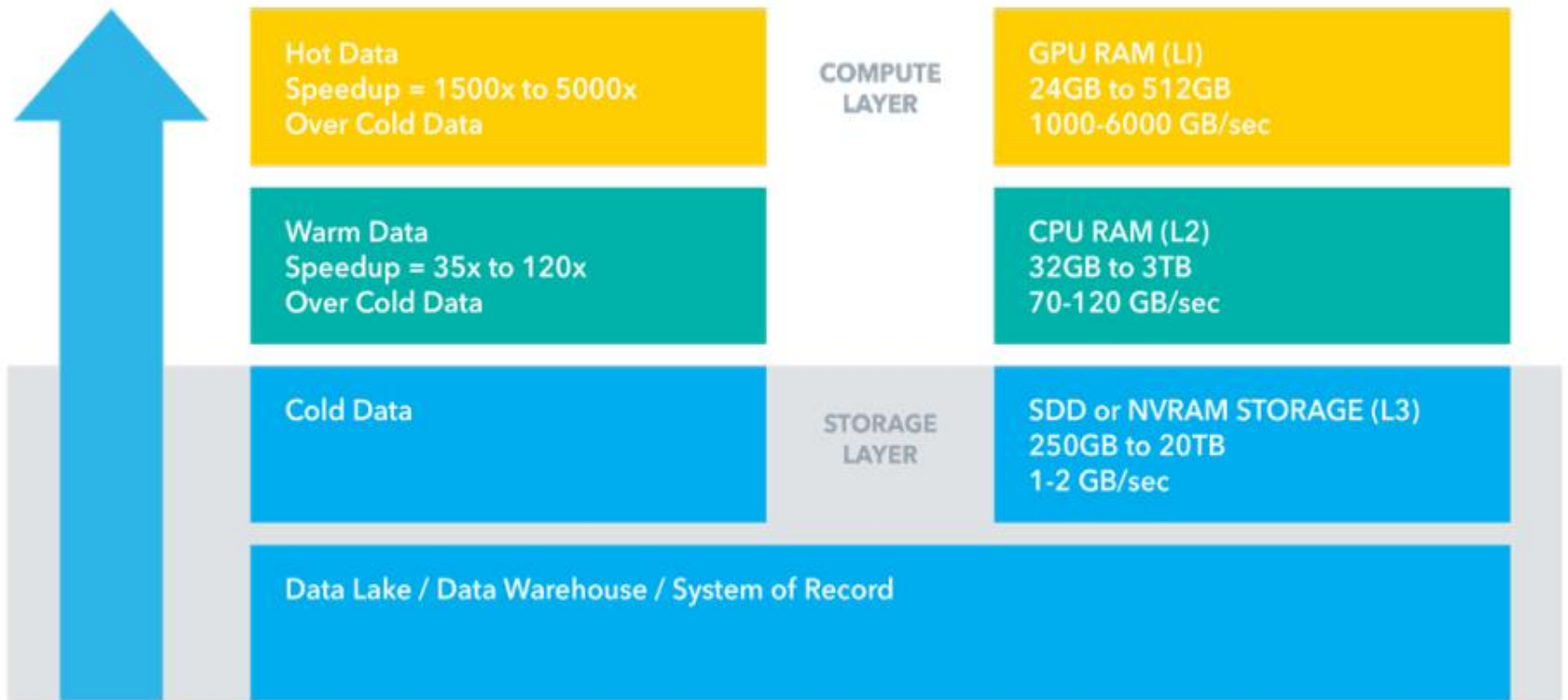


OmniSciDB

OmniSciDB является основой платформы OmniSci. OmniSciDB основана на SQL, является реляционной, столбчатой и специально разработана для использования массового параллелизма современного оборудования CPU и GPU. OmniSciDB может запрашивать до миллиардов строк за миллисекунды и таким образом удовлетворяет требованиям обработки больших, высокоскоростных данных.

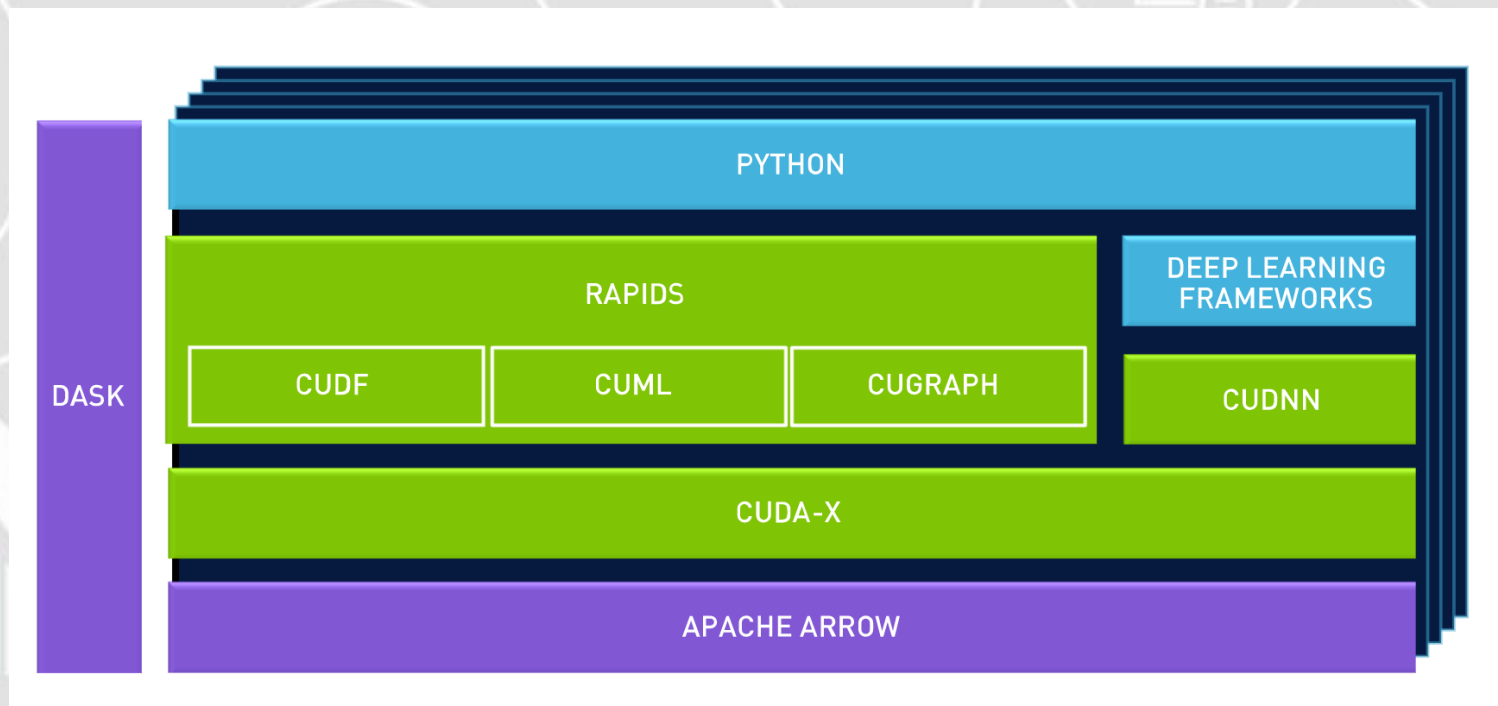


Структура и характеристики OmniSciDB



Rapids

RAPIDS включает набор открытых библиотек для анализа, машинного обучения и визуализации данных с GPU-ускорением. Эта платформа разрабатывалась инженерами NVIDIA в сотрудничестве с ключевыми разработчиками открытого ПО.



Высокоуровневые API

Python

Dask Multi-GPU ML

Scikit-Learn-Like

Dask-CUML

Data Distribution

CuML

CUDA/C++

ML Algorithms

ML Primitives

libcuml

Model Parallelism

Multi-Node / Multi-GPU Communications

Host 1

GPU1

GPU3

GPU2

GPU4

Host 2

GPU1

GPU3

GPU2

GPU4

Arrow

Apache Arrow - это многоязычная платформа для разработки данных в памяти. Она определяет стандартизированный независимый от языка формат колоночной памяти для плоских и иерархических данных, организованных для эффективных аналитических операций на современном оборудовании.

	session_id	timestamp	source_ip
Row 1	1331246660	3/8/2012 2:44PM	99.155.155.225
Row 2	1331246351	3/8/2012 2:38PM	65.87.165.114
Row 3	1331244570	3/8/2012 2:09PM	71.10.106.181
Row 4	1331261196	3/8/2012 6:46PM	76.102.156.138

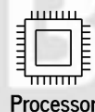
Tradicional Memory Buffer

Row 1	1331246660	3/8/2012 2:44PM	99.155.155.225
Row 2	1331246351	3/8/2012 2:38PM	65.87.165.114
Row 3	1331244570	3/8/2012 2:09PM	71.10.106.181
Row 4	1331261196	3/8/2012 6:46PM	76.102.156.138

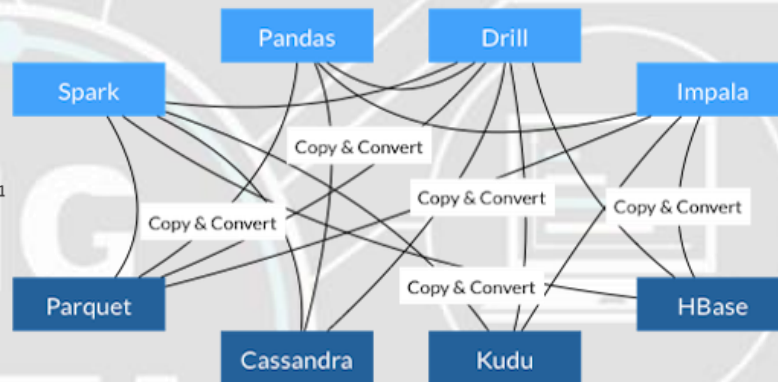
Arrow Memory Buffer

session_id	1331246660	1331246351	1331244570	1331261196
timestamp	3/8/2012 2:44PM	3/8/2012 2:38PM	3/8/2012 2:09PM	3/8/2012 6:46PM
source_ip	99.155.155.225	65.87.165.114	71.10.106.181	76.102.156.138

SELECT * FROM clickstream
WHERE session_id = 1331246351



Processor



CUDF

CUDF-это библиотека фреймов данных Python GPU (построенная на столбчатом формате памяти Apache Arrow) для загрузки, объединения, агрегирования, фильтрации и иного манипулирования табличными данными.



**BIG
DATA**

Типы данных CUDF

Kind of Data	Data Type	Scalar	String Aliases
Integer		<code>np.int8</code> , <code>np.int16</code> , <code>np.int32</code> , <code>np.int64</code> , <code>np.uint8</code> , <code>np.uint16</code> , <code>np.uint32</code> , <code>np.uint64</code>	<code>'int8'</code> , <code>'int16'</code> , <code>'int32'</code> , <code>'int64'</code> , <code>'uint8'</code> , <code>'uint16'</code> , <code>'uint32'</code> , <code>'uint64'</code>
Float		<code>np.float32</code> , <code>np.float64</code>	<code>'float32'</code> , <code>'float64'</code>
Strings		<code>str</code>	<code>'string'</code> , <code>'object'</code>
Datetime		<code>np.datetime64</code>	<code>'datetime64[s]'</code> , <code>'datetime64[ms]'</code> , <code>'datetime64[us]'</code> , <code>'datetime64[ns]'</code>
Timedelta (duration type)		<code>np.timedelta64</code>	<code>'timedelta64[s]'</code> , <code>'timedelta64[ms]'</code> , <code>'timedelta64[us]'</code> , <code>'timedelta64[ns]'</code>
Categorical	CategoricalDtype	(none)	<code>'category'</code>
Boolean		<code>np.bool</code>	<code>'bool'</code>

Примеры работы в CUDF

```
import cudf, io, requests
from io import StringIO

url = "https://github.com/plotly/datasets/raw/master/tips.csv"
content = requests.get(url).content.decode('utf-8')

tips_df = cudf.read_csv(StringIO(content))
tips_df['tip_percentage'] = tips_df['tip'] / tips_df['total_bill'] * 100

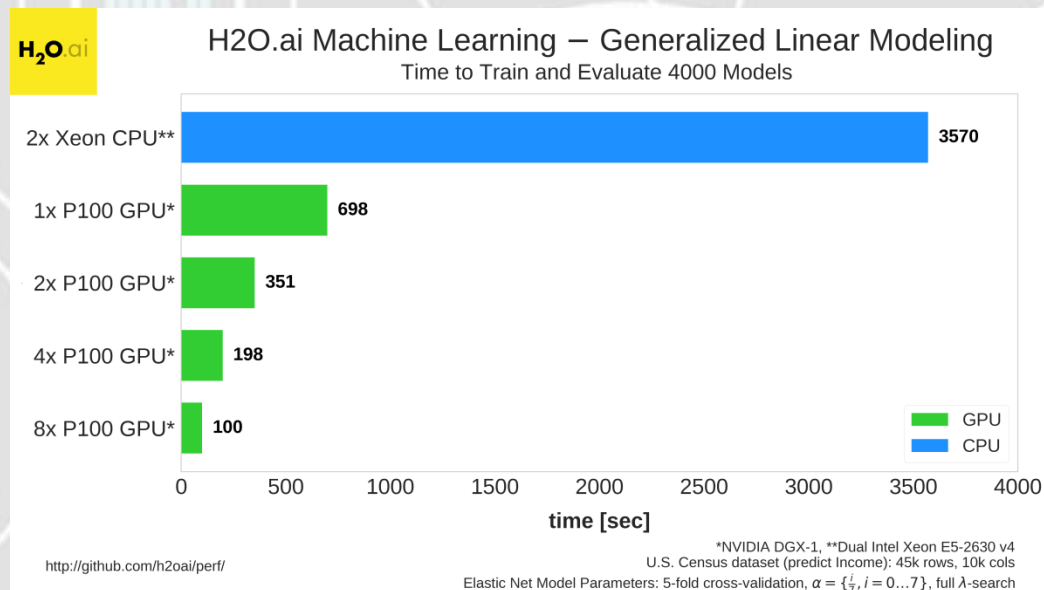
# display average tip by dining party size
print(tips_df.groupby('size').tip_percentage.mean())
```

Output:

```
size
1    21.729201548727808
2    16.571919173482897
3    15.215685473711837
4    14.594900639351332
5    14.149548965142023
6    15.622920072028379
Name: tip_percentage, dtype: float64
```

H2O4GPU

H2O4GPU-это пакет машинного обучения с открытым исходным кодом, ускоренный GPU, с API на Python и R, который позволяет любому пользователю использовать преимущества графических процессоров для создания передовых моделей машинного обучения. Существует множество популярных алгоритмов, включая машины градиентного бустинга (GBM), обобщенные линейные модели (GLM) и кластеризацию К-средних.



<https://github.com/h2oai/h2o4gpu>

H2O4GPU

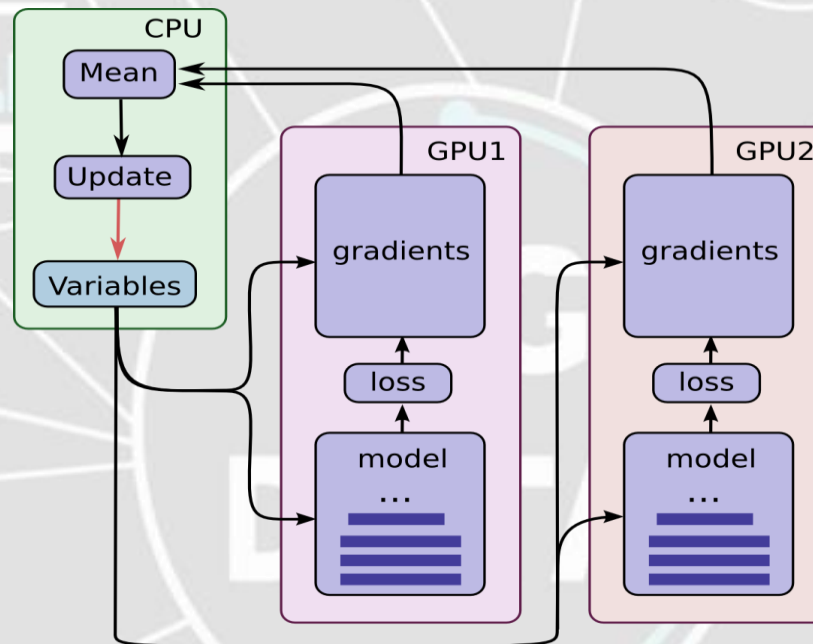
```
import h2o4gpu  
import numpy as np
```

```
X = np.array([[1.,1.], [1.,4.], [1.,0.]])  
model = h2o4gpu.KMeans(n_clusters=2,random_state=1234).fit(X)  
model.cluster_centers_
```

**BIG
DATA**

Tensorflow-GPU

TensorFlow-это библиотека программного обеспечения с открытым исходным кодом для высокопроизводительных численных вычислений. Ее гибкая архитектура позволяет легко развертывать вычисления на различных платформах (процессоры, графические процессоры, TPU), а также на настольных компьютерах, кластерах серверов и мобильных и периферийных устройствах.



Tensorflow-GPU

```
1
2 !pip -q install tensorflow==1.13.1
3 !pip install tensorflow-gpu==1.13.1
4 !pip -q install opencv-python
5 !pip -q install keras==2.3.0
6 !pip -q install ImageAi==2.1.5
7 !wget https://github.com/OlafenwaMoses/ImageAI/releases/download/1.0/resnet50_coco_best_v2.0.1.h5
8 !python -c 'from tensorflow.python.client import device_lib;print(device_lib.list_local_devices());'
```

```
1 from imageai.Prediction.Custom import ModelTraining
2 model_trainer = ModelTraining()
3 model_trainer.setModelTypeAsResNet()
4 model_trainer.setDataDirectory("/content/odezhda")
5 model_trainer.trainModel(num_objects=6, num_experiments=100, enhance_data=True, batch_size=32, show_network_summary=True)
```