



Microsoft ML for Apache Spark

Unifying Machine Learning Ecosystems at Massive Scales

Mark Hamilton

Microsoft, MIT

marhamil@microsoft.com

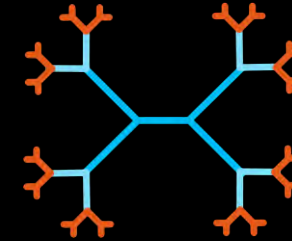
Overview



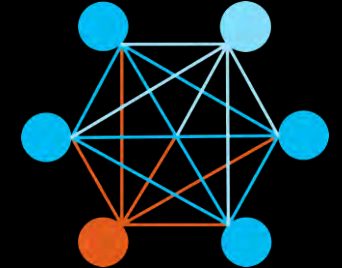
- ▶ Background
 - ▶ Spark + SparkML
 - ▶ MMLSpark
- ▶ Unifying ML Ecosystems
 - ▶ LightGBM, CNTK, Vowpal Wabbit
 - ▶ Multilingual Bindings
- ▶ Microservice Orchestration
 - ▶ Cognitive Services on Spark
- ▶ Model Deployment with Spark Serving
- ▶ Use Cases
 - ▶ The Snow Leopard Trust



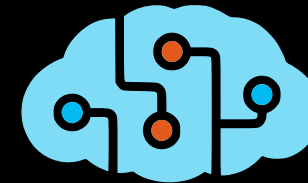
Vowpal Wabbit



LightGBM



CNTK



Cognitive Services



Kubernetes

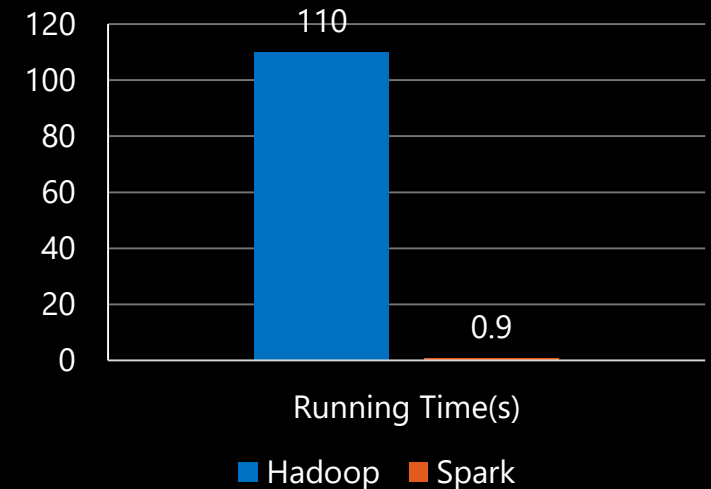
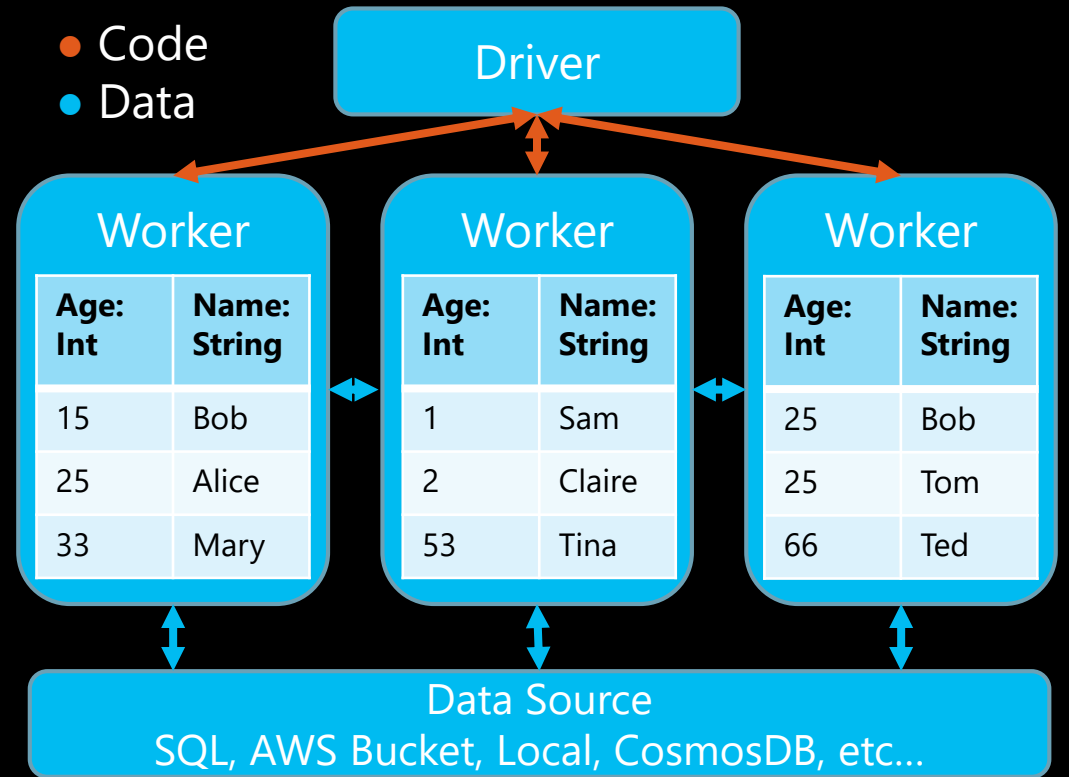


Snow
Leopard
Trust





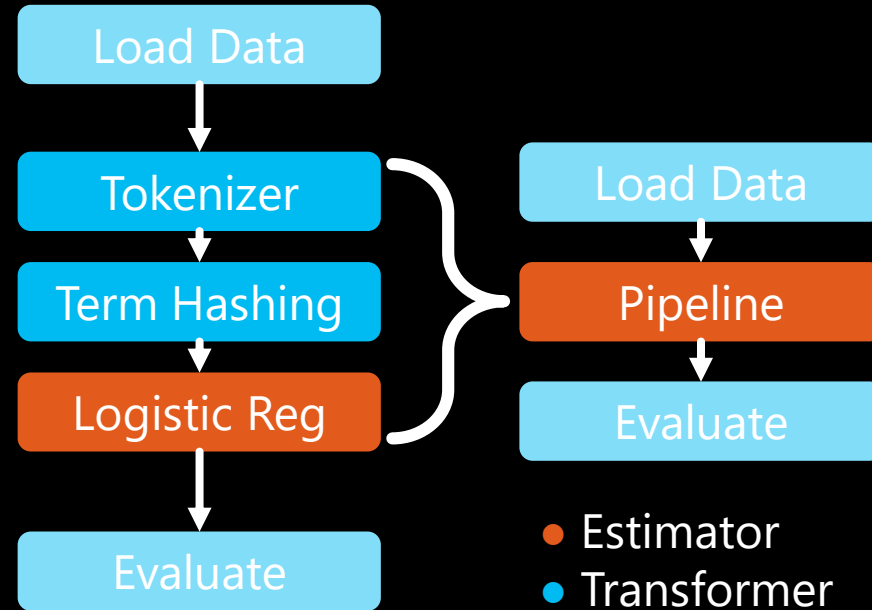
- ▶ A **fault-tolerant distributed** computing framework
- ▶ Map Reduce + SQL
- ▶ Whole program optimization + query pushdown
- ▶ Elastic
- ▶ Scala, Python, R, Java, Julia
- ▶ ML, Graph Processing, Streaming





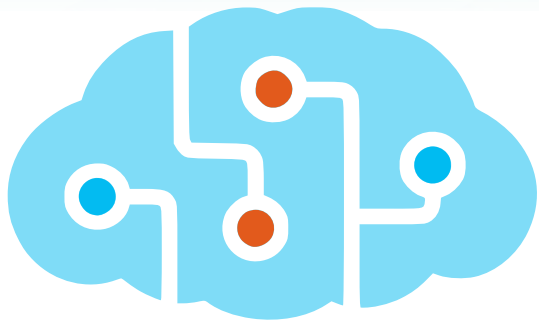
- ▶ High level library for distributed machine learning
- ▶ More general than SciKit-Learn
- ▶ All models have a uniform interface
 - ▶ Can compose models into complex pipelines
 - ▶ Can save, load, and transport models

```
data = spark.read.csv("hdfs://...")
train, test = data.randomSplit([.5, .5])
model = LogisticRegression().fit(train)
predictions = model.transform(test)
```



Microsoft Machine Learning for Apache Spark v0.18

Microsoft's Open Source Contributions to Apache Spark



Distributed
Machine Learning



Fast Model
Deployment



Microservice
Orchestration



Multilingual Binding
Generation

www.aka.ms/spark

 [Azure/mmlspark](https://github.com/Azure/mmlspark)

Unifying Machine Learning Ecosystems

▶ Goals

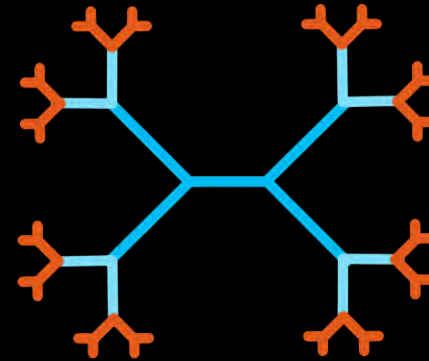
- ▶ Same API
- ▶ Composable
- ▶ Batch, Streaming, Serving
- ▶ Elastically Distributed
- ▶ Fault Tolerant
- ▶ Multi-Language
- ▶ Data Source Agnostic



Markus Cozowicz
marcozo@microsoft.com
Data Scientist



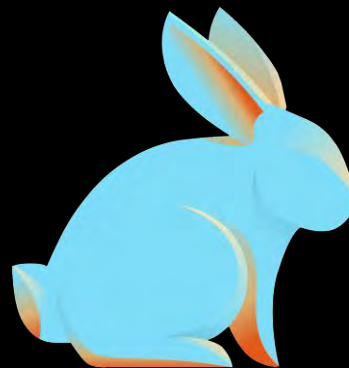
Image Processing
with Open CV



Gradient Boosting
with LightGBM



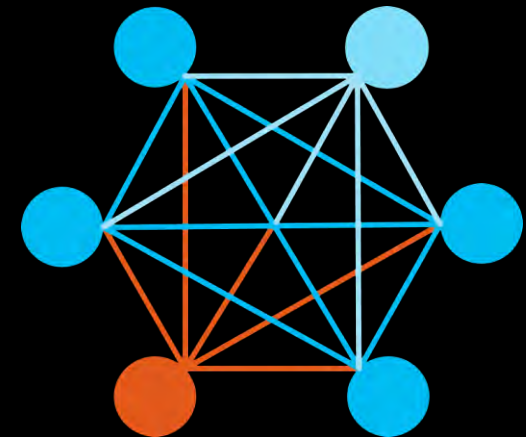
Deep Learning
Pipelines (Databricks)



**Text Analytics with
Vowpal Wabbit**



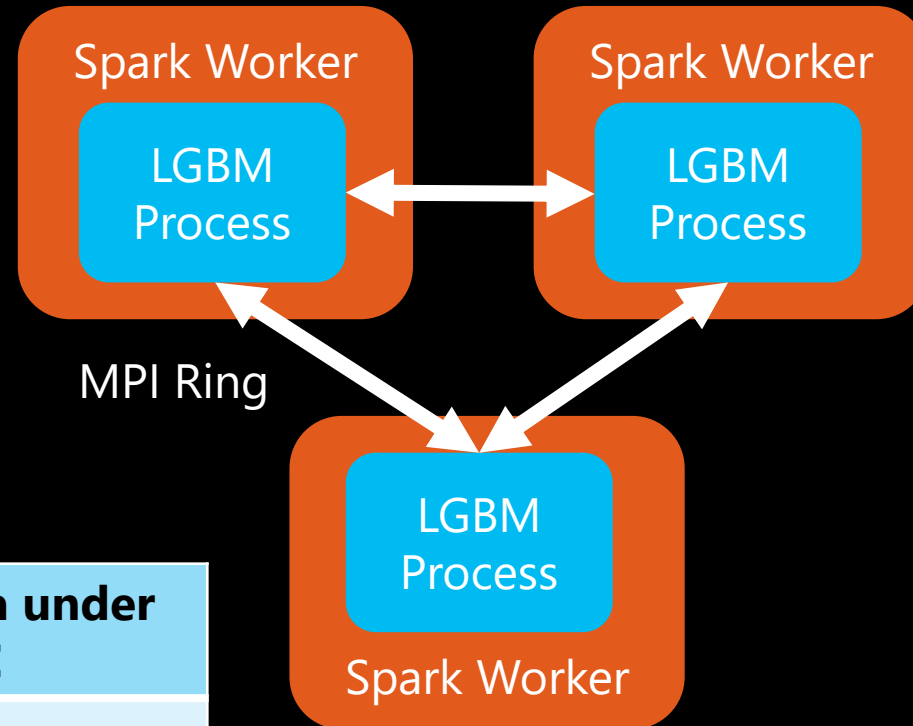
Distributed Model
Interpretability with LIME



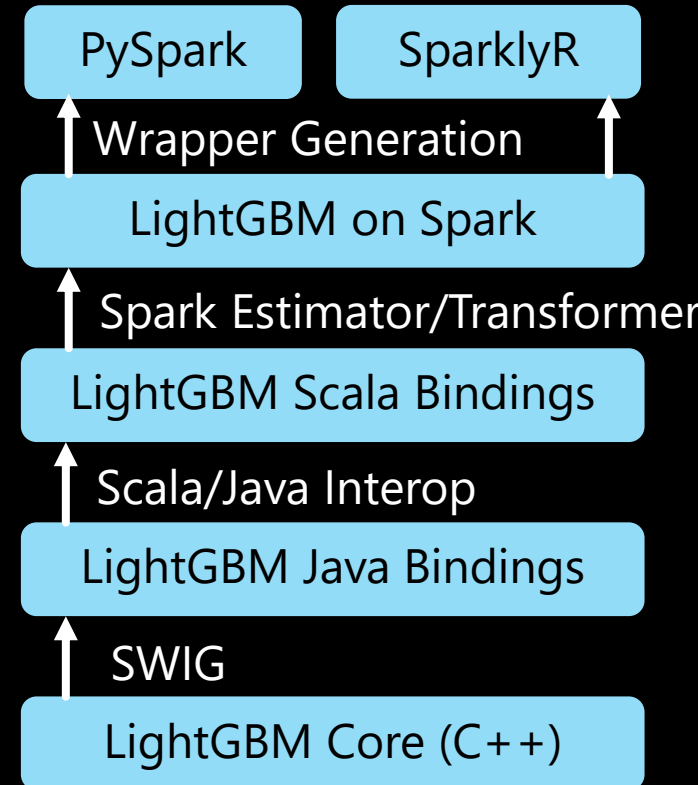
Deep Learning
with CNTK

Example Backend: LightGBM on Spark

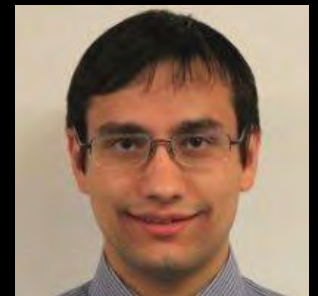
- ▶ Barrier Execution for Synchronizing Workers
- ▶ Fast Socket/MPI communication
- ▶ mapPartitions for Transformer



Framework	Time(s)	Area under ROC
XGBoost	52.60	.808
SparkML GBT	82.78	.788
LightGBM	45.39	.812

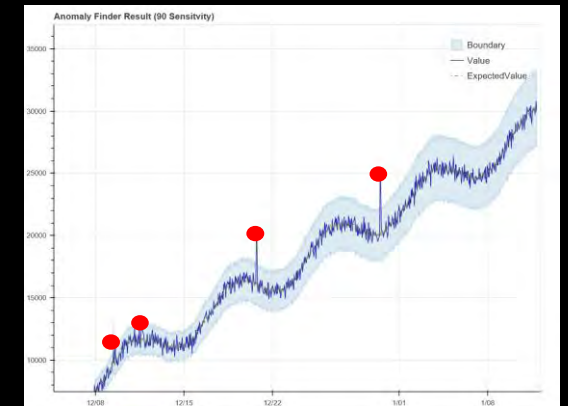
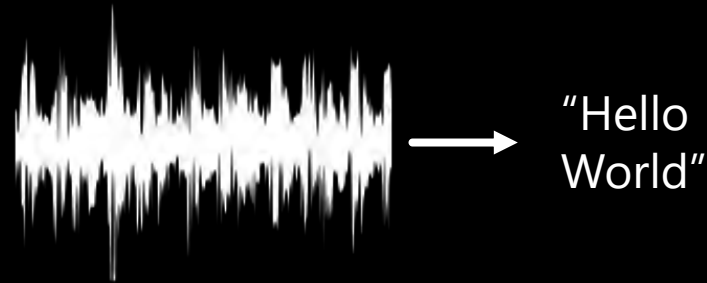
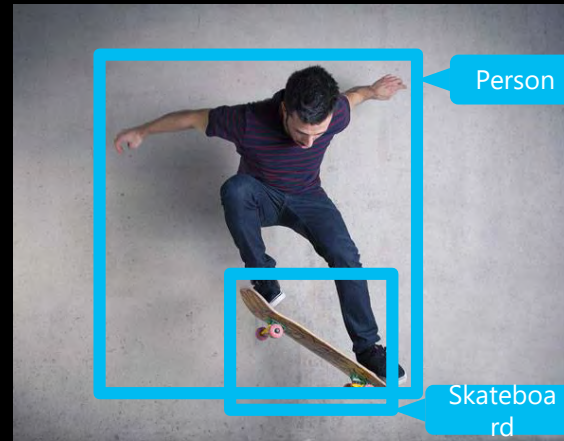


Ilya Matiach, ilmat@microsoft.com
Developer, Azure ML



Cognitive Services

- ▶ High quality pre-built intelligent services
- ▶ No time intensive model training or deployment
- ▶ Leverage Microsoft Research and Azure ML
- ▶ **Available as Docker Containers**



I had a wonderful trip to Seattle last week and even visited the Space Needle 2 times!

Place Time Range

En-US 84% positive



Vision

- Object, scene, and activity detection
- Face recognition and identification
- Celebrity and landmark recognition
- Emotion recognition
- Text and handwriting recognition (OCR)
- Customizable image recognition
- Video metadata, audio, and keyframe extraction and analysis
- Explicit or offensive content moderation



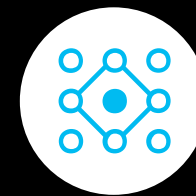
Speech

- Speech transcription (speech-to-text)
- Custom speech models for unique vocabularies or complex environment
- Text-to-speech
- Custom Voice
- Real-time speech translation
- Customizable speech transcription and translation
- Speaker identification and verification



Language

- Language detection
- Named entity recognition
- Key phrase extraction
- Text sentiment analysis
- Multilingual and contextual spell checking
- Explicit or offensive text content moderation
- PII detection for text moderation
- Text translation
- Customizable text translation
- Contextual language understanding



Decision

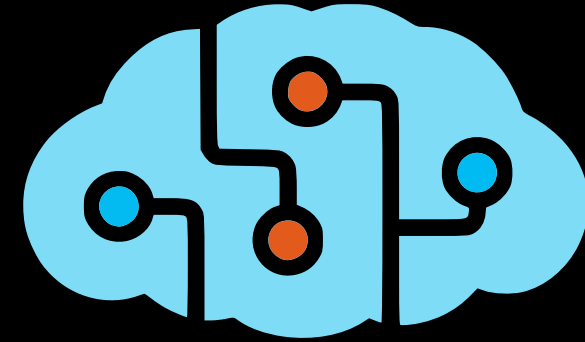
- Q&A extraction from unstructured text
- Knowledge base creation from collections of Q&As
- Semantic matching for knowledge bases
- Customizable content personalization learning



Search

- Ad-free web, news, image, and video search results
- Trends for video, news
- Image identification, classification and knowledge extraction
- Identification of similar images and products
- Named entity recognition and classification
- Knowledge acquisition for named entities
- Search query autosuggest
- Ad-free custom search engine creation

Azure Cognitive Services on Spark



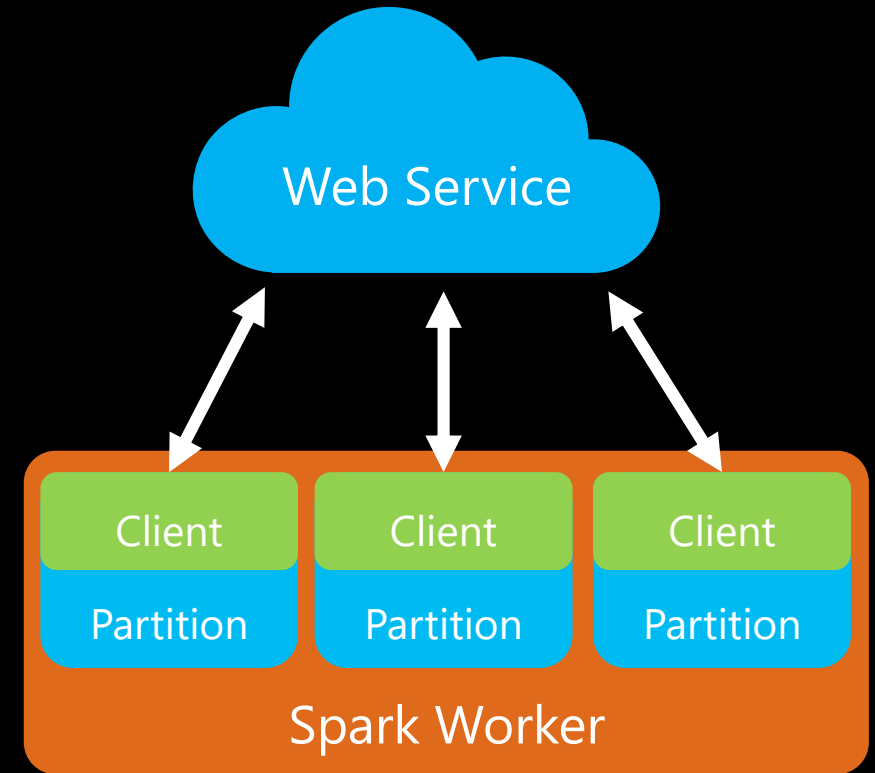
- ▶ Easy to use integration between Spark and the Azure Cognitive Services
- ▶ Composable and pipelinable with all other SparkML models!
- ▶ Exponential Backoffs, Backpressure, Batching, Async Parallelism
- ▶ Fully Fluent API

```
val df = new TextSentiment()  
    .setTextCol("text")  
    .setOutputCol("sentiment")  
    .transform(inputs)
```

Features	Time (s)	Errors #
None	30.8	18993
EBO+BP	1163.0	0
EBO+BP+B	57.1	0
EBO+BP+B+P	49.7	0

HTTP on Spark

- ▶ Full Integration between HTTP Protocol and Spark SQL
- ▶ Spark as a Microservice Orchestrator
- ▶ Spark + X
- ▶ Support for all Spark Languages



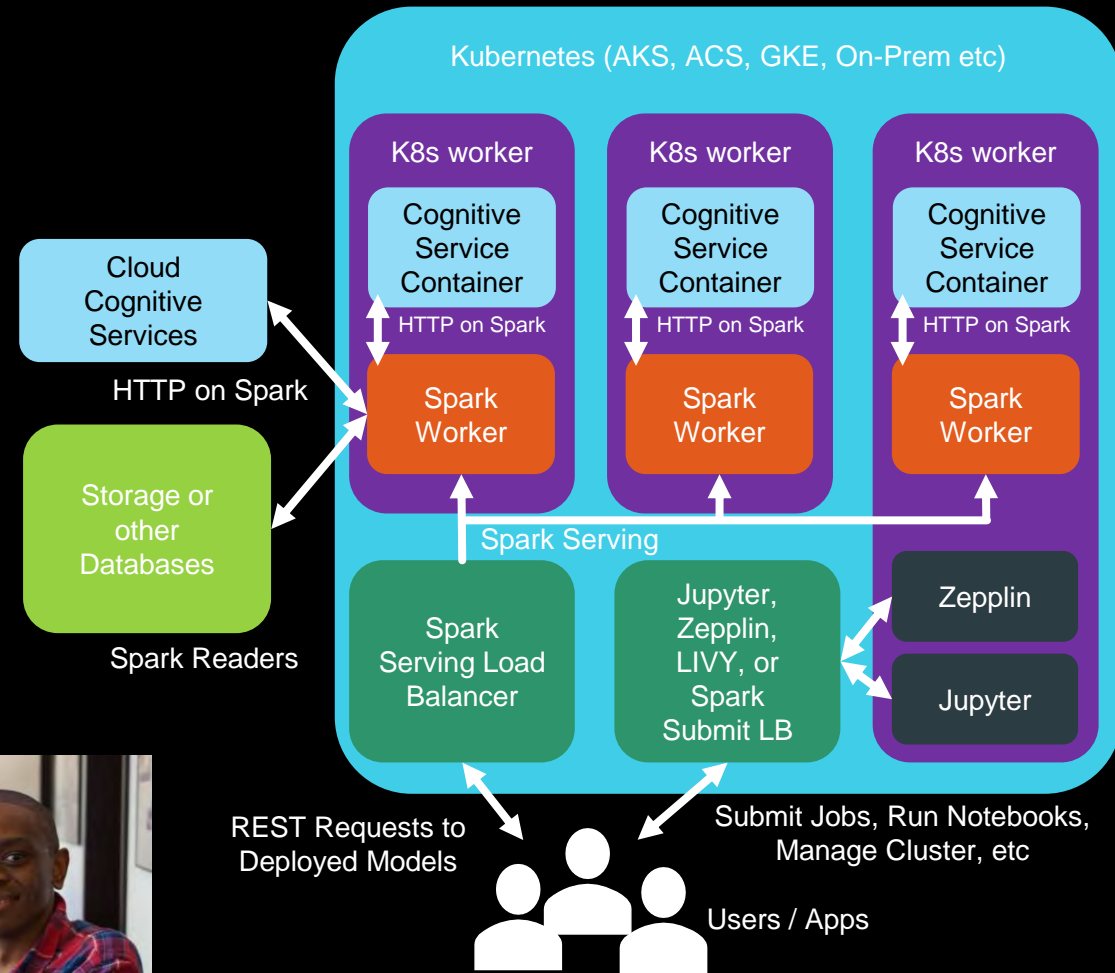
```
df = SimpleHTTPTransformer()  
    .setInputParser(JSONInputParser())  
    .setOutputParser(JSONOutputParser()  
        .setDataType(schema))  
    .setOutputCol("results")  
    .setUrl(...)
```

Deploying on Kubernetes

- ▶ Works on any k8s cluster
- ▶ Helm: Package Manager for Kubernetes

```
helm repo add microsoft \
  https://microsoft.github.io/charts/repo
helm update
```

```
helm install microsoft/spark --version 1.0.0
```



Dalitso Banda, dbanda@microsoft.com
Microsoft AI Development Acceleration Program

Model Deployment with Spark Serving



- ▶ Sub-millisecond RESTful Model Deployment on Spark Clusters

Batch API:

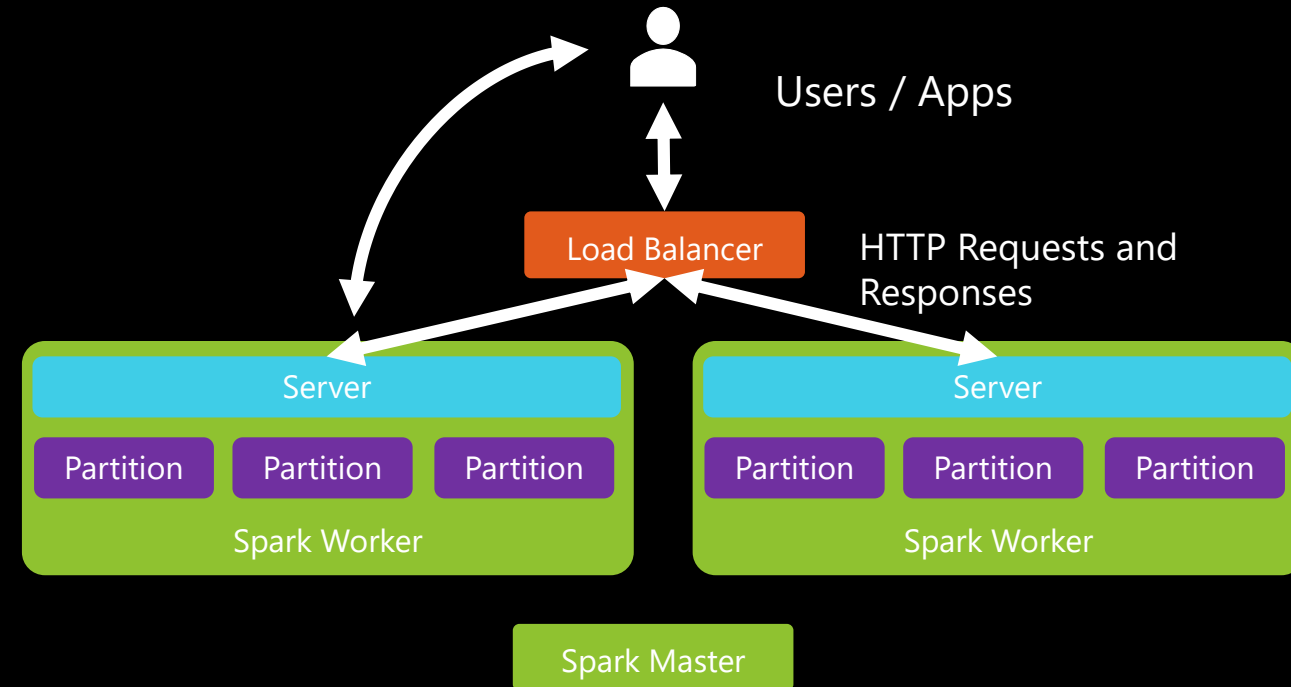
```
spark.read.parquet.load(...)  
  .select(...)
```

Streaming API:

```
spark.readStream.kafka.load(...)  
  .select(...)
```

Serving API:

```
spark.readStream.server("0.0.0.0", 5000).load(...)  
  .select(...)
```



AI for Earth



Snow
Leopard
Trust

Endangered Status Matters

BBC Sign in News Sport Weather Shop Earth Travel

NEWS


Home Video World US & Canada UK Business Tech Science Stories Enter

Asia China India

Snow leopard no longer 'endangered'

14 September 2017

f t m Share



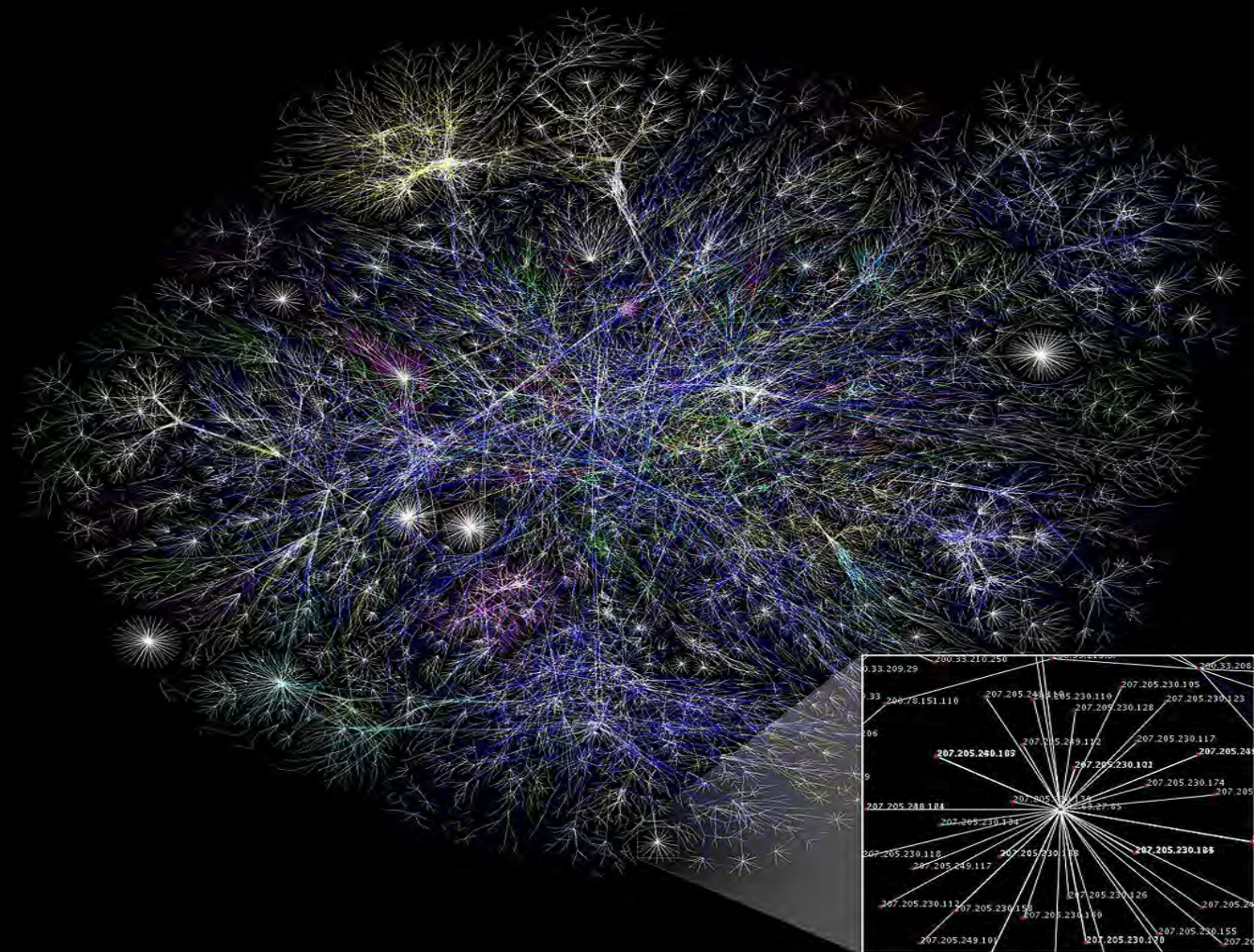
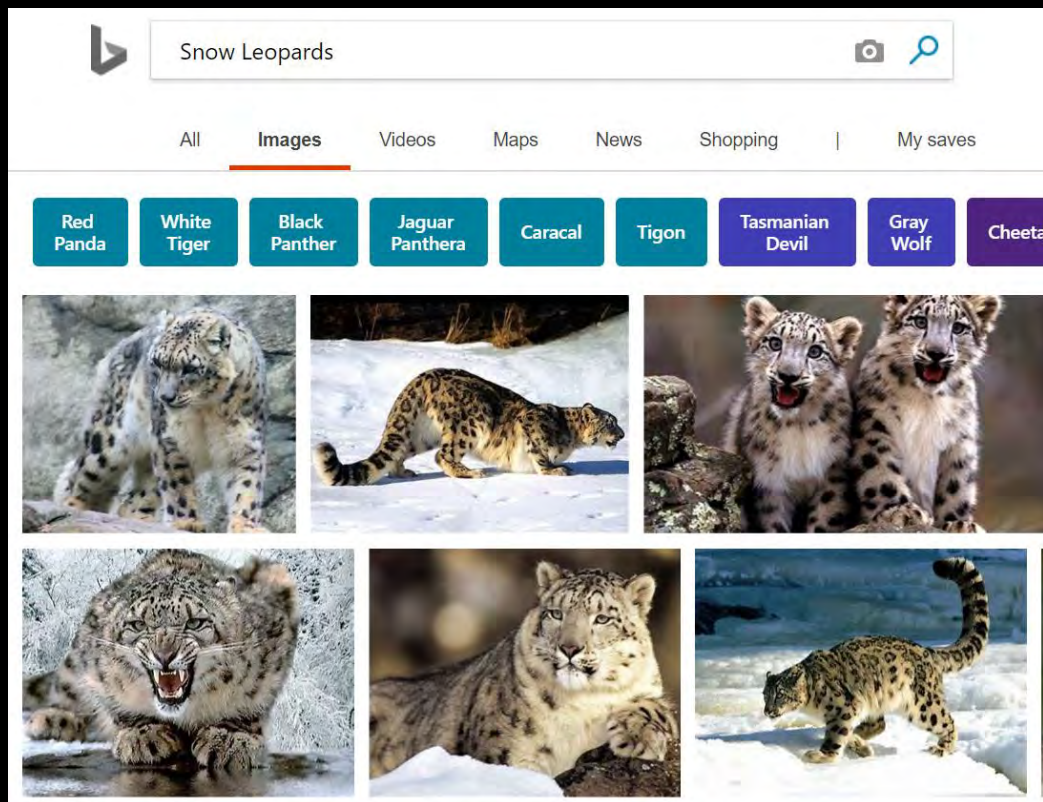
Statement on IUCN Red List Status Change of the Snow Leopard

The Snow Leopard Trust, one the leading conservation organizations working to protect this cat, opposes the IUCN's decision to change the snow leopard's Red List status from 'Endangered' to 'Vulnerable'.

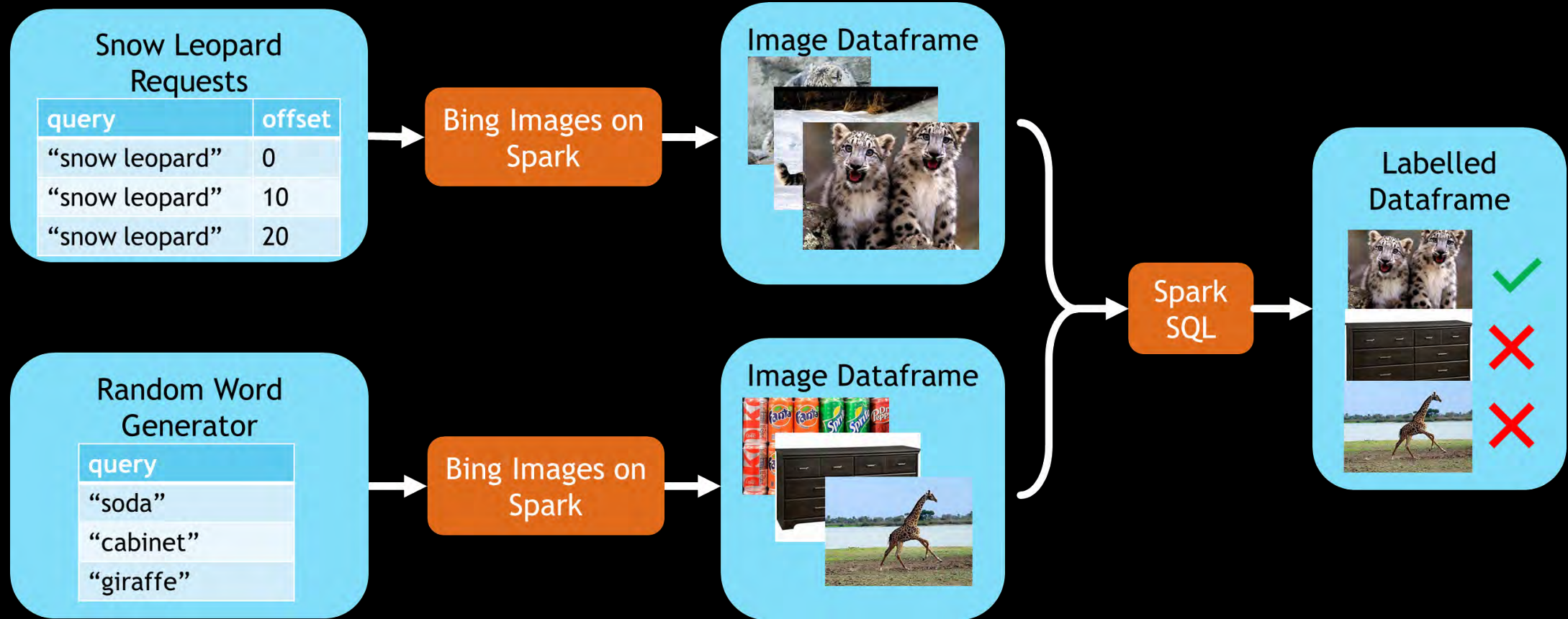
Remote Camera Trapping



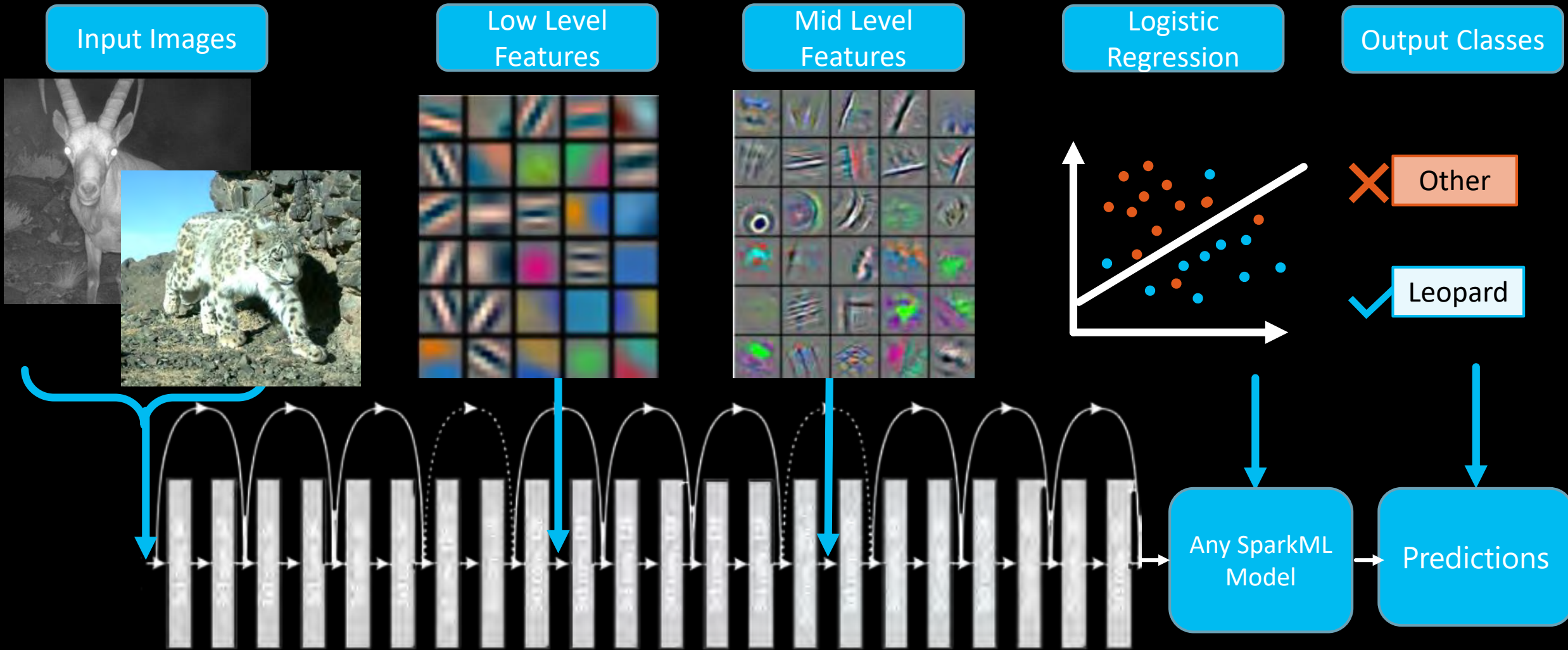
Creating a labelled Training Dataset



Creating a labelled Training Dataset



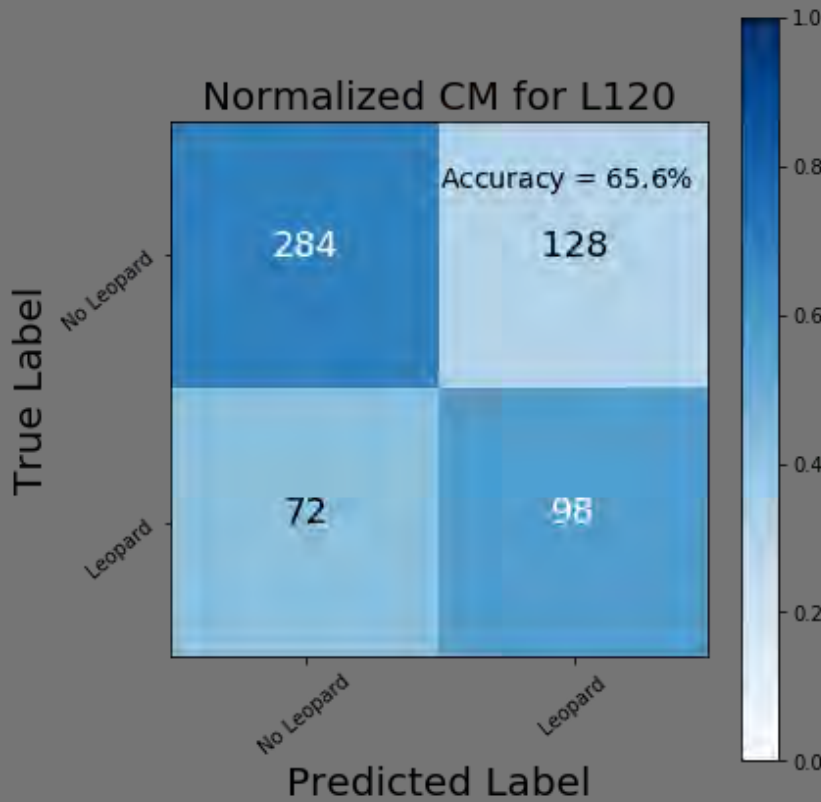
Transfer Learning with ResNet 50



Filters from Zeiler + Fergus 2013

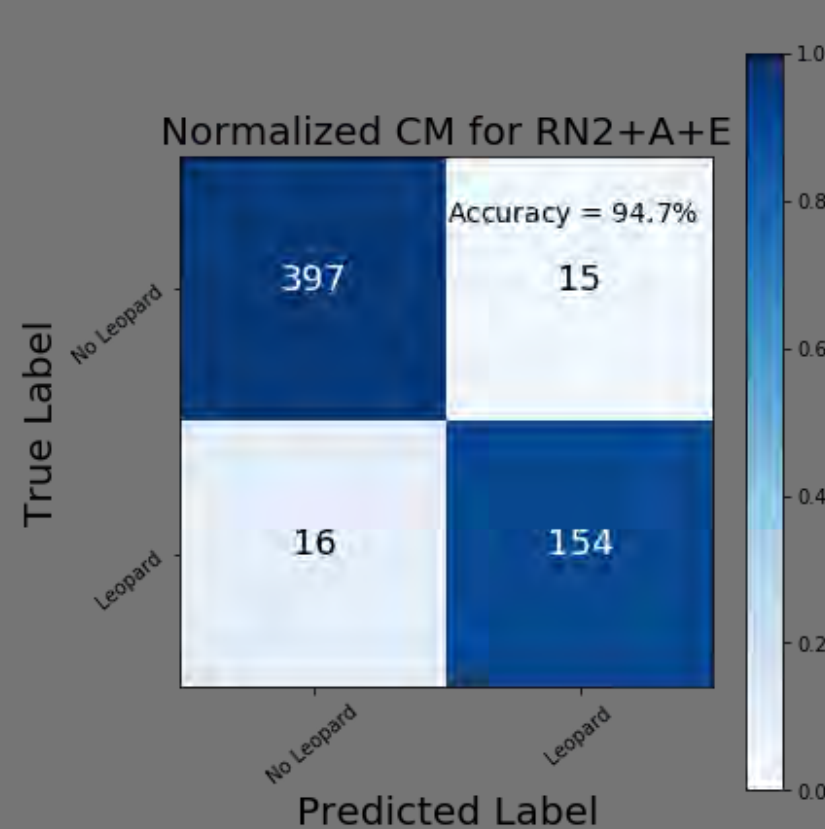
Performance

Without Deep Featurization

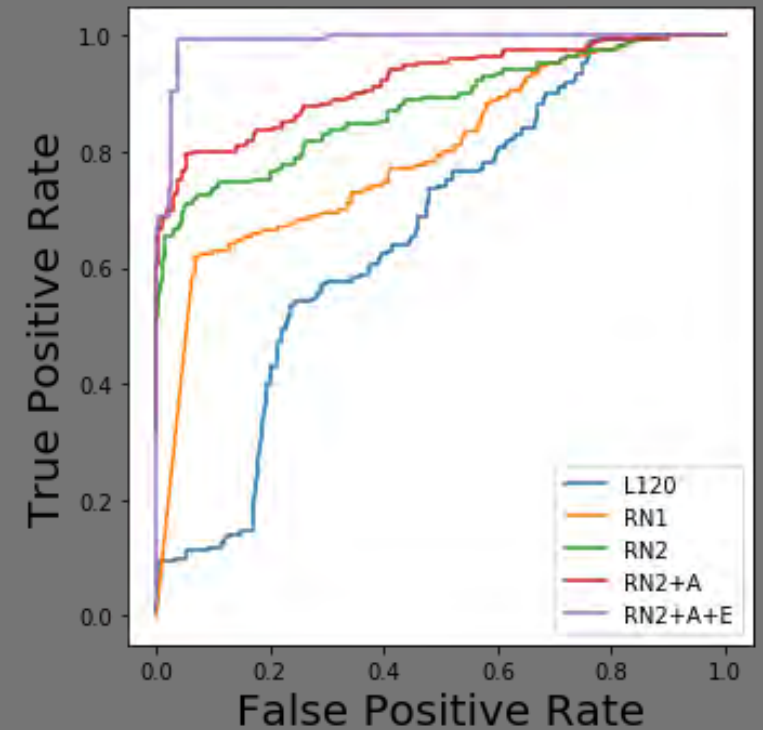


Accuracy 65.6%

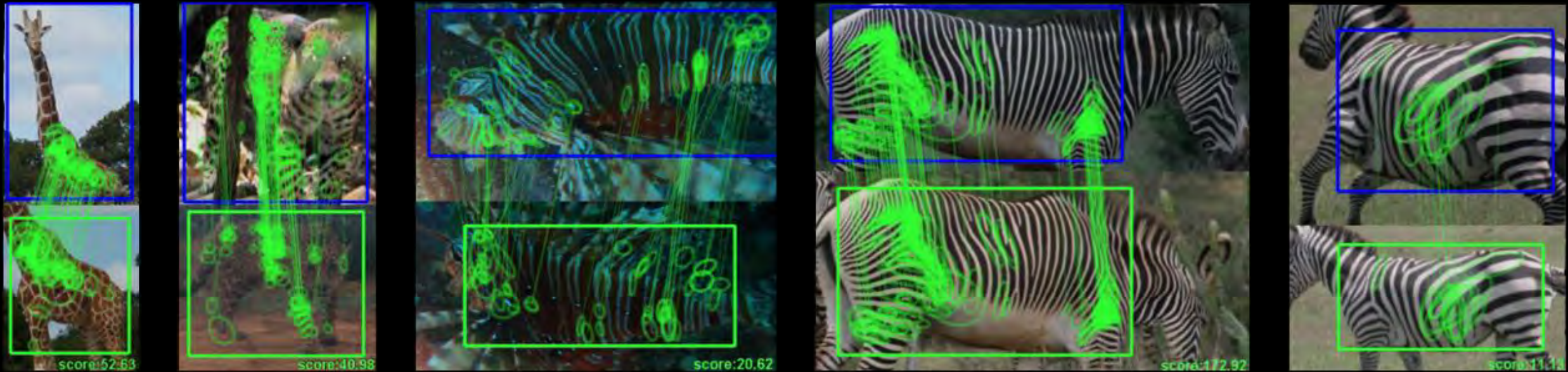
With Deep Featurization, Augmentation, and Temporal Ensembling



Accuracy 94.7%



Goal: Identify Individual Leopards






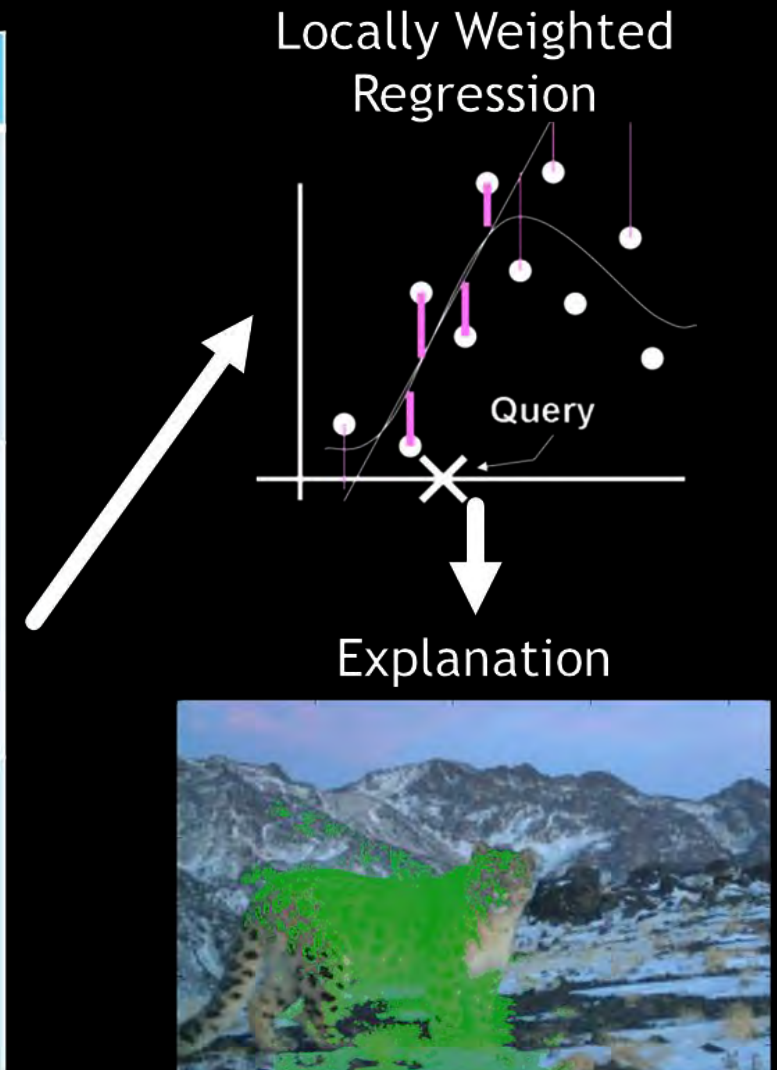
Source: HotSpotter - Patterned Species Instance Recognition

Automating Detection with LIME on Spark



$P(\text{leopard}) = .88$

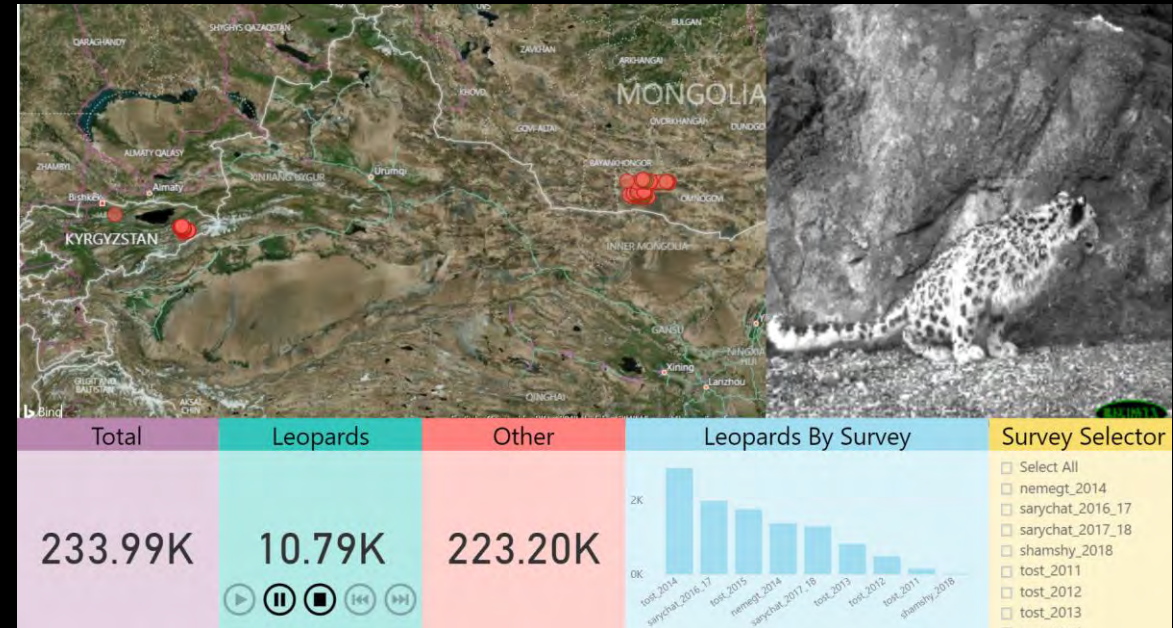
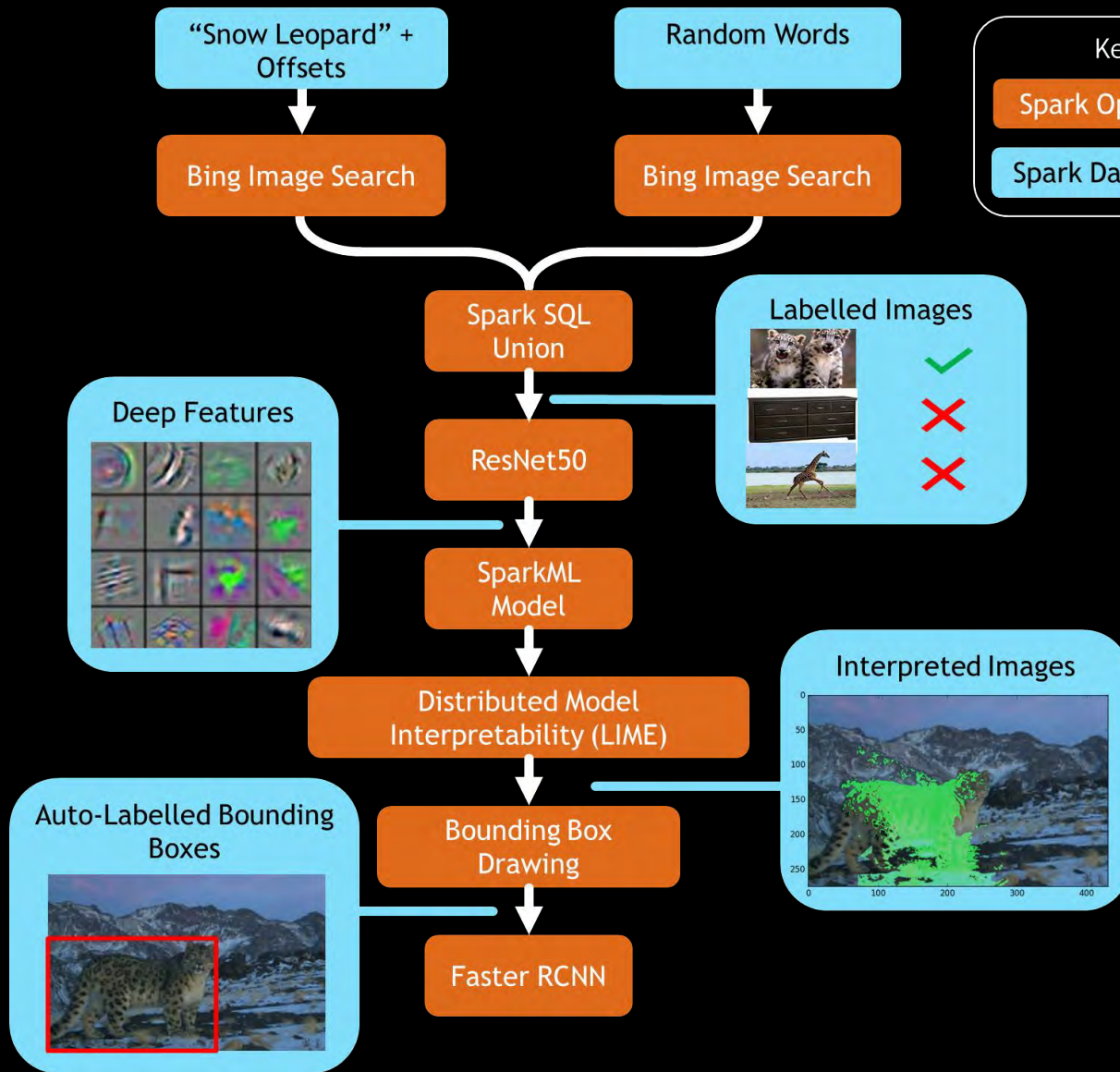
Perturbed Photos	$P(\text{leopard})$
	.92
	.56
	.003



LIME on Spark



End to End Architecture



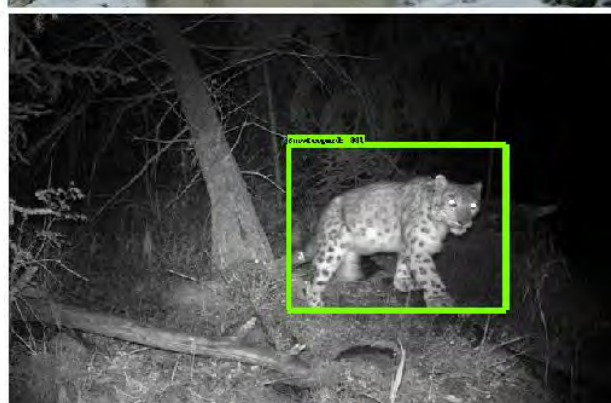
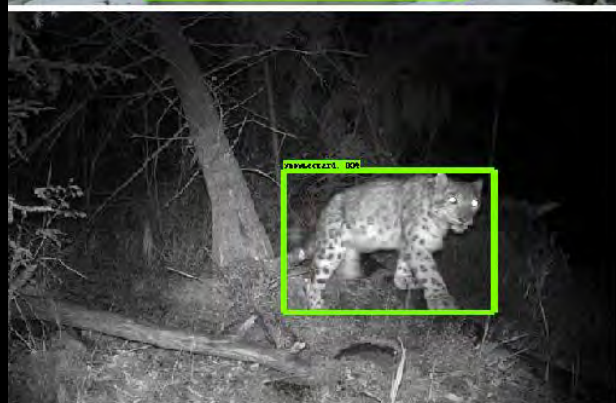
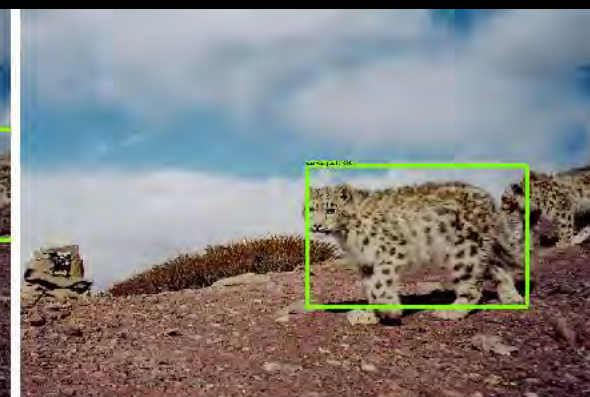
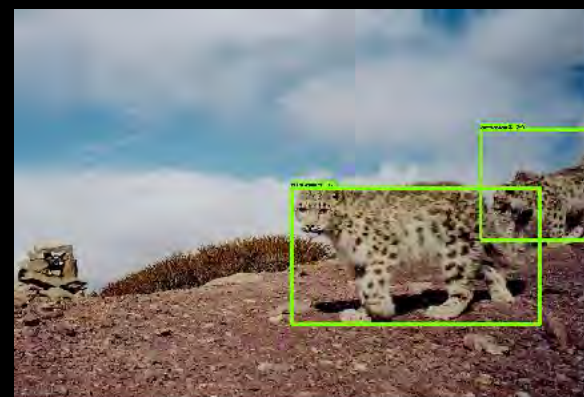
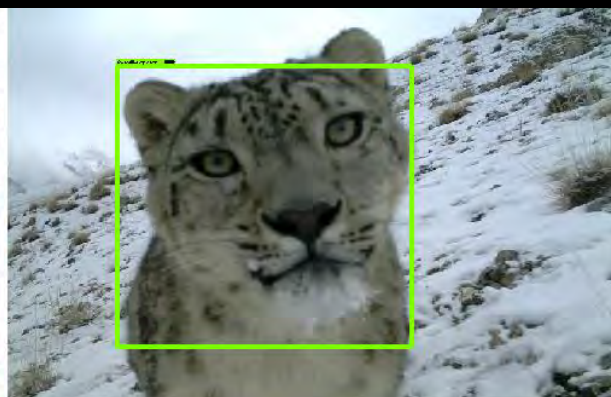
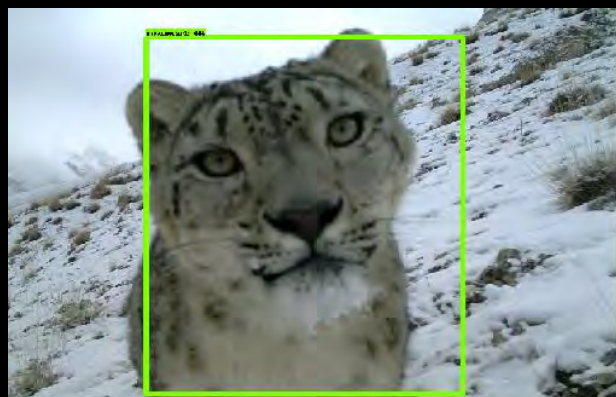
Results

Human Labels

Unsupervised
FRCNN Outputs

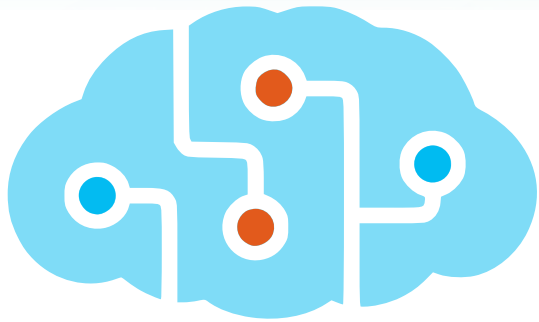
Human Labels

Unsupervised
FRCNN Outputs



Microsoft Machine Learning for Apache Spark v0.18

**Microsoft's Open Source
Contributions to Apache Spark**



Distributed
Machine Learning



Fast Model
Deployment



Microservice
Orchestration



Multilingual Binding
Generation


www.aka.ms/spark

 [Azure/mmlspark](https://github.com/Azure/mmlspark)

Thanks to

- ▶ You all!
- ▶ **Ilya Matiach: LightGBM on Spark**
- ▶ **Markus Cozowicz: VW on Spark**
- ▶ Sudarshan Raghunathan, Christina Lee, Daniel Ciborowski, Eli Barzilay, Tong Wen, Pablo Castro, Chris Hoder, Ryan Gaspar, Henrik Neilsen, Andrew Schonhoffer, Joseph Sirosh
- ▶ Microsoft NERD Garage Team + MIT Externship Program
- ▶ Snow Leopard Trust: Koustubh Sharma, Rhetick Sengupta, Jeff Brown, Michael Despines
- ▶ Microsoft Development Acceleration Team:
 - ▶ Dalitso Banda, Casey Hong, Karthik Rajendran, Manon Knoertzer, Tayo Amuneke, Alejandro Buendia
- ▶ Azure CAT, AzureML, and Azure Search Teams

Get in Touch

- ▶ Support: mmlspark-support@microsoft.com
- ▶ Me: marhamil@microsoft.com
- ▶ Github  : Azure/mmlspark
- ▶ Website: www.aka.ms/spark
- ▶ Paper: www.aka.ms/spark-paper
- ▶ Contributions Welcome!
- ▶ Check out our MSR Podcast on Oct 2

Backup Slides

AI for Cultural Institutions



Celebrating 2 years of Open Access at The MET

- ▶ In 2016 The MET Released 400k images under open access
- ▶ This past winter the MET released a new subject-keyword dataset of image annotations
- ▶ MIT, The MET, and Microsoft participated in a 3-day hackathon to create intelligent experiences using the new collection



Open Access



Goals:

Create new
works of art

Use new work
to explore
existing art

Explore further
with intelligent
search

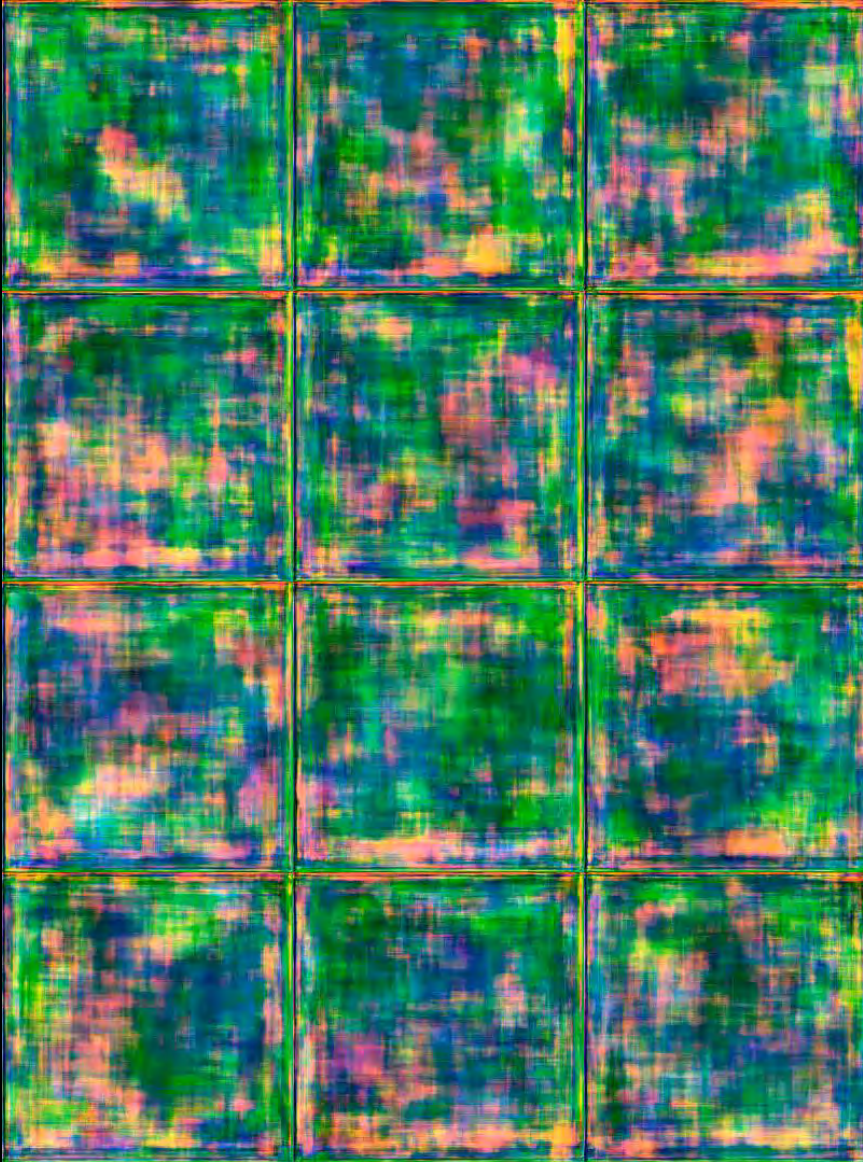
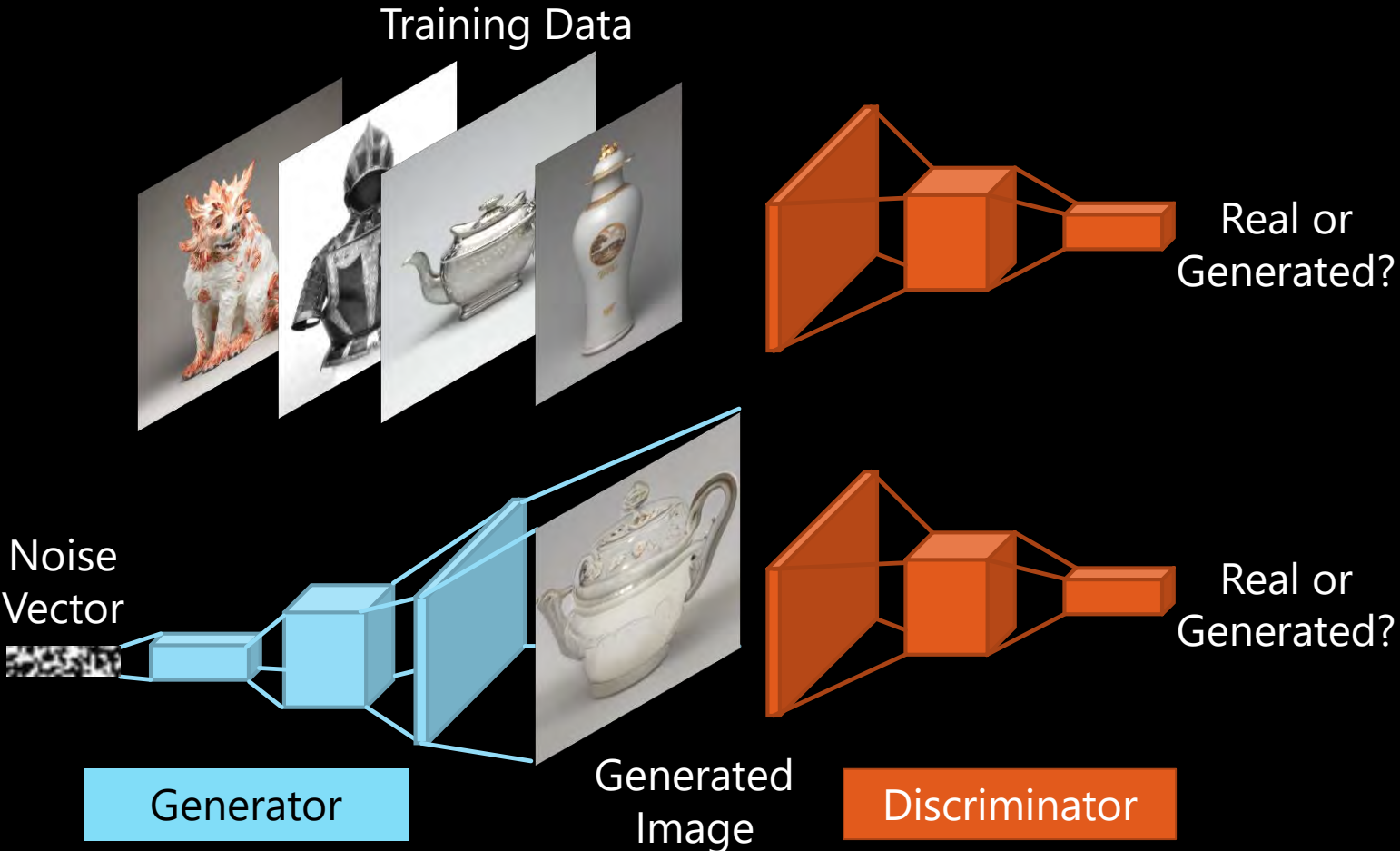
Needed Technologies:

Generative
Adversarial
Networks

Reverse image
search

Elasticsearch
with Cognitive
Services

Generative Adversarial Art



Custom Reverse Image Search

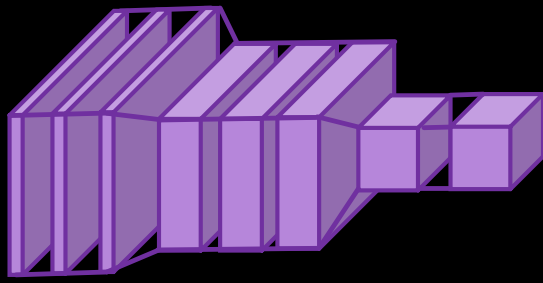
Query Image

ResNet Featurizer

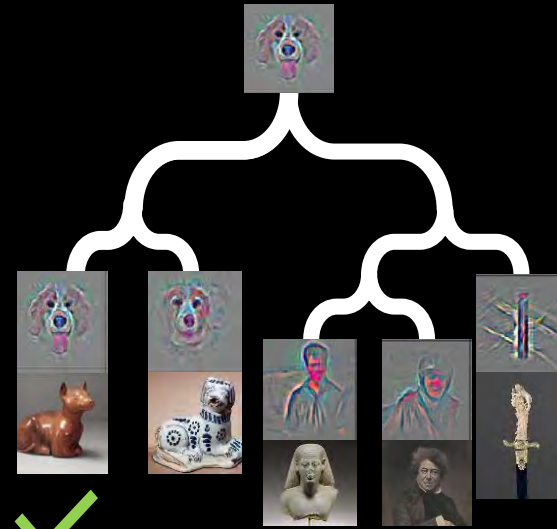
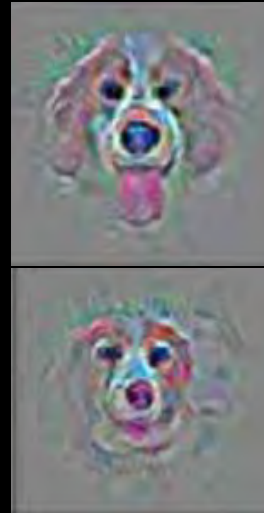
Deep Features

Fast Nearest Neighbor Lookup

Closest Match



MMLSpark



SparkML LSH or Annoy



Example Nearest Neighbors



Intelligent Search Index

- ▶ Pipe images through Computer Vision API to annotate image for searching
- ▶ Stream images and intelligent annotations to Azure Search



Query Image:



Describe Image Output:

A picture containing a person

Deep Feature Nearest Neighbors:



A picture containing a glass, cup




A fish swimming underwater




End Application: Gen Studio

Generated Image




Click or drag to explore the space between them











EXPLORE SIMILAR

SAVE IMAGE



tiger

- Watercolor
- Drawings
- British
- Lustreware
- Ceramics-Pottery
- Watercolor on heavy wove paper
- Watercolor and gouache
- Iran
- Ceramic, paint
- Ceramics-Sculpture
- Pen and brush and iron gall ink
- French

 <p>Tiger Drawings and Prints</p>	 <p>Tiger European Sculpture and Decorative Arts</p>	 <p>Tiger Drawings and Prints</p>	 <p>Royal Tiger Drawings and Prints</p>
			

AI for Accessibility



Seeing AI



Currency Identification

1928-1934 Series* - 1928 to 1935	1928-1934 Series* - 1928 to 1935	1928-1934 Series* - 1928 to 1935	1928-1934 Series* - 1928 to 1935	1928-1934 Series* - 1928 to 1935	1928-1934 Series* - 1928 to 1935
1935-1956 Series* - 1935 to 1956	1935-1956 Series* - 1935 to 1956	1935-1956 Series* - 1935 to 1956	1935-1956 Series* - 1935 to 1956	1935-1956 Series* - 1935 to 1956	1935-1956 Series* - 1935 to 1956
1957-1963 Series - 1957 to Date	1957-1963 Series - 1957 to Date	1957-1963 Series - 1957 to Date	1957-1963 Series - 1957 to Date	1957-1963 Series - 1957 to Date	1957-1963 Series - 1957 to Date
1963-1969 Series - 1963 to August 16, 1969 (discontinued)	1963-1969 Series - 1963 to August 16, 1969 (discontinued)	1963-1969 Series - 1963 to August 16, 1969 (discontinued)	1963-1969 Series - 1963 to August 16, 1969 (discontinued)	1963-1969 Series - 1963 to August 16, 1969 (discontinued)	1963-1969 Series - 1963 to August 16, 1969 (discontinued)
1969-2009 Series - 1969 to Date	1969-2009 Series - 1969 to Date	1969-2009 Series - 1969 to Date	1969-2009 Series - 1969 to Date	1969-2009 Series - 1969 to Date	1969-2009 Series - 1969 to Date
2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date
1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date
1928-1934 Series - 1928 to July 14, 1935 (discontinued)	1928-1934 Series - 1928 to July 14, 1935 (discontinued)	1928-1934 Series - 1928 to July 14, 1935 (discontinued)	1928-1934 Series - 1928 to July 14, 1935 (discontinued)	1928-1934 Series - 1928 to July 14, 1935 (discontinued)	1928-1934 Series - 1928 to July 14, 1935 (discontinued)
1935-1956 Series - 1935 to July 14, 1956 (discontinued)	1935-1956 Series - 1935 to July 14, 1956 (discontinued)	1935-1956 Series - 1935 to July 14, 1956 (discontinued)	1935-1956 Series - 1935 to July 14, 1956 (discontinued)	1935-1956 Series - 1935 to July 14, 1956 (discontinued)	1935-1956 Series - 1935 to July 14, 1956 (discontinued)
1957-1963 Series - 1957 to July 14, 1969 (discontinued)	1957-1963 Series - 1957 to July 14, 1969 (discontinued)	1957-1963 Series - 1957 to July 14, 1969 (discontinued)	1957-1963 Series - 1957 to July 14, 1969 (discontinued)	1957-1963 Series - 1957 to July 14, 1969 (discontinued)	1957-1963 Series - 1957 to July 14, 1969 (discontinued)
1963-1969 Series - 1963 to July 14, 1969 (discontinued)	1963-1969 Series - 1963 to July 14, 1969 (discontinued)	1963-1969 Series - 1963 to July 14, 1969 (discontinued)	1963-1969 Series - 1963 to July 14, 1969 (discontinued)	1963-1969 Series - 1963 to July 14, 1969 (discontinued)	1963-1969 Series - 1963 to July 14, 1969 (discontinued)
1969-2009 Series - 1969 to July 14, 1969 (discontinued)	1969-2009 Series - 1969 to July 14, 1969 (discontinued)	1969-2009 Series - 1969 to July 14, 1969 (discontinued)	1969-2009 Series - 1969 to July 14, 1969 (discontinued)	1969-2009 Series - 1969 to July 14, 1969 (discontinued)	1969-2009 Series - 1969 to July 14, 1969 (discontinued)
2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date	2009 Series - 2009 to Date
1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date	1976 Series - 1976 to Date

MODERN UNITED STATES CURRENCY



REPRODUCED FROM THE FEDERAL RESERVE'S U.S. CURRENCY...
 PART OF THE FEDERAL RESERVE'S U.S. CURRENCY...
 © 2013 U.S. Treasury Department. All rights reserved.
 The Federal Reserve Bank of New York is not responsible for the accuracy of the information provided on this page.
 The information on this page is for informational purposes only and is not intended to be used as a substitute for professional advice.
 The information on this page is for informational purposes only and is not intended to be used as a substitute for professional advice.
 The information on this page is for informational purposes only and is not intended to be used as a substitute for professional advice.

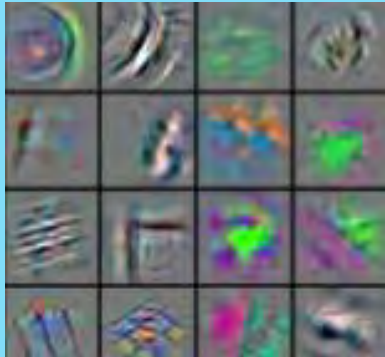
A Familiar Architecture...

Labelled Images



\$1
\$5
\$20

Deep Features



Queries
1 Dollar
5 Dollars
10 Dollars
20 Dollars

Bing Image Search

Union + Distinct

MobileNet

Logistic Regression

Azure Machine Learning + Spark Serving



Prep Data

Train

Deploy