

Deep Video Analytics

A data-centric approach to Computer Vision

Akshay Bhat
Cornell Tech, Cornell University.

Quick overview of Computer Vision over last two decades

<http://www.computervisionblog.com/2015/01/from-feature-descriptors-to-deep.html>

SIFT, Graph Cuts



HOG, DPM



Deep Learning



?

Caltech 101, Matlab, OpenCV



VOC, Imagenet, Caffe, Theano



?

Developments over last 5 years

High quality libraries

- OpenCV
- ROS
- Caffe
- Theano
- Torch
- Tensor Flow
- CNTK
- MXNET
- PyTorch
- deeplearn.js

Developments over last 5 years

Pre-trained models

- Imagenet classification
 - Inception
 - Resnet
 - VGG
- Detection models
 - R-CNN
 - YOLO
 - SSD
- Face detection / recognition
 - Face-MTCNN
 - Facenet
- Semantic Segmentation models
 - Multipathnet
 - FCN
- Audio embedding models
 - Soundnet

Developments over last 5 years

A deluge of datasets!

- Open Images
- Yahoo Flickr Creative Com. 100M
- MSCOCO
- ViCom
- Visual Genome
- YouTube-BoundingBoxes / 8M
- AMOS
- imSitu, Charades by AllenAI
- KITTI /Toronto City
- Udacity car dataset
- Caltech, INRIA, ETH Pedestrians
- Stanford Drone Dataset
- Uber text
- THUMOS

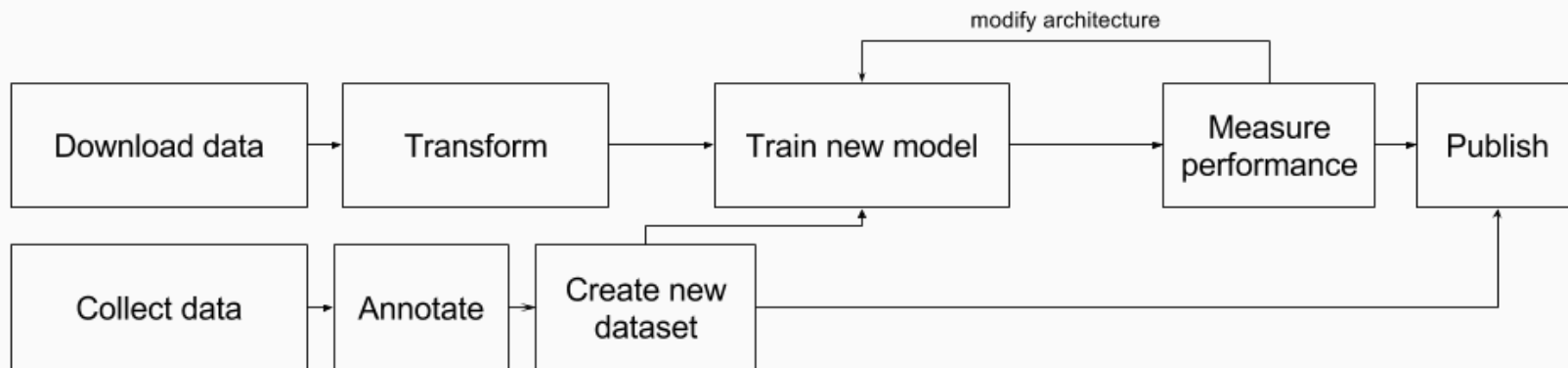
Number of datasets \approx Number of research groups

With each dataset having its own JSON or XML format, incompatible with all others.

What is hidden in plain sight?

Model-centric approach

Libraries & frameworks are designed with **goal of training and evaluation of models for individual tasks.**

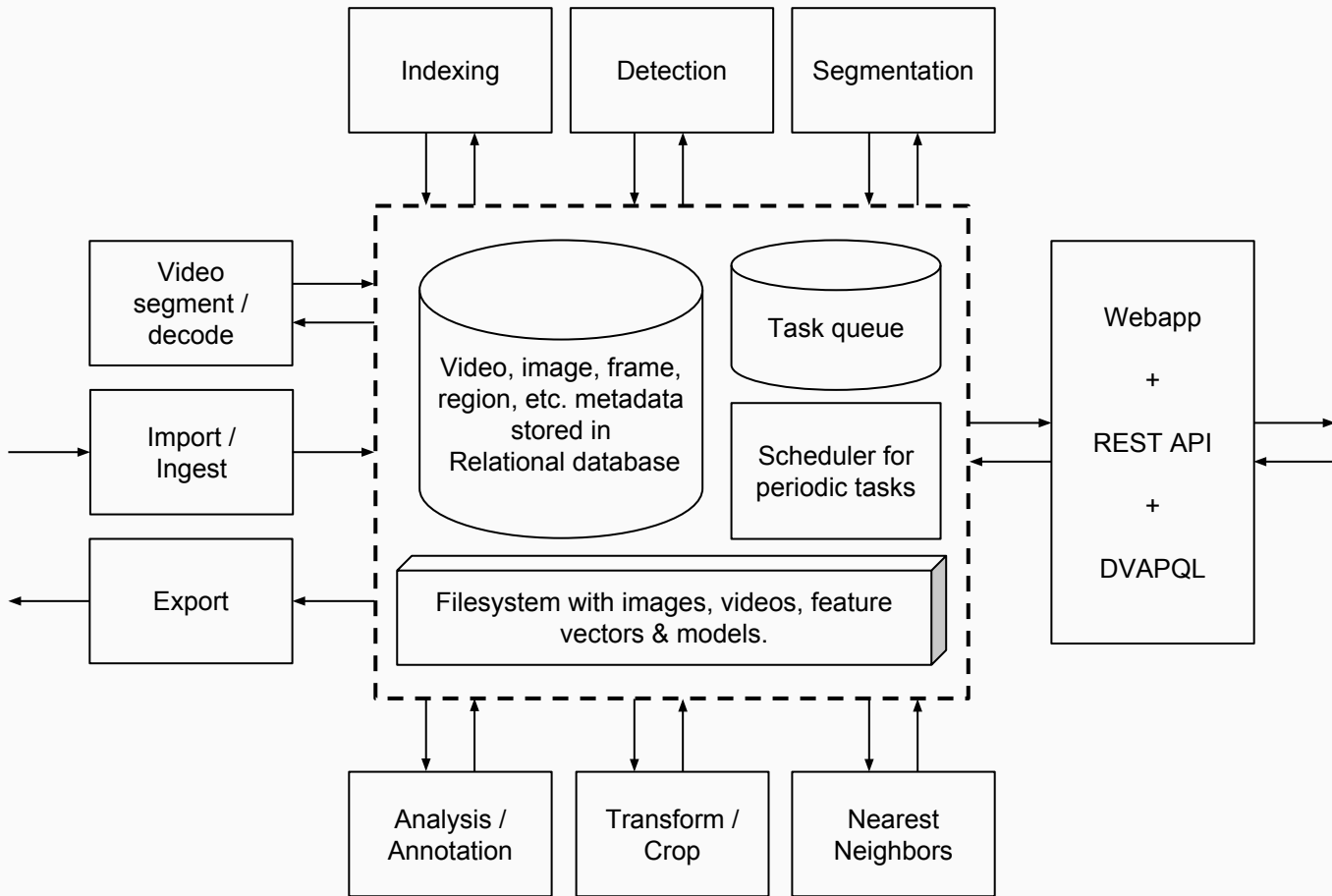


Unsuitable for building systems that learn in interactive manner, or leverage data from multiple sources or combine multiple tasks.

We need a data-centric approach that allows us to combine

- Models for multiple tasks
- Data from multiple sources
- User Interaction / interface

Model-centric to Data-centric



A Relational Model of Data for Large Shared Data Banks. By Edgar F. Codd

Can we develop an equivalent of relational model for visual data?

Relational data : Postgres, MYSQL, SQLite

::

Text, HTML : Lucene/Solr, Elasticsearch

::

Videos & Images : _____

Previous attempts: LIRE project

- LIRE: Lucene Image Retrieval
 - <http://www.lire-project.net/>
- Developed pre-Deep Learning
- Functionality limited to computing & storing feature vectors such as Color Layout, Edge Histogram, etc.

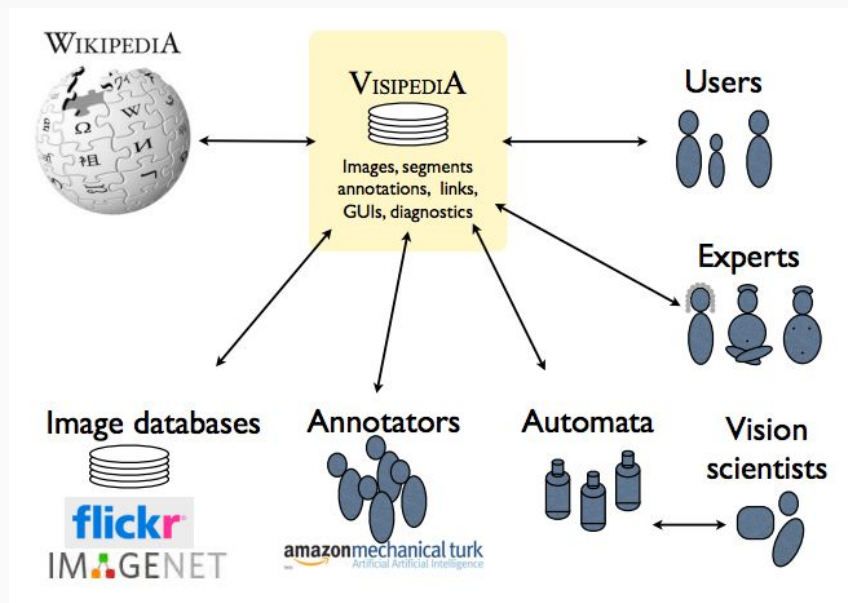
Previous attempts: CloudCV

- Large Scale Distributed Computer Vision as a Cloud Service
- Support for OpenCV, Graphlab, Cafe
- Image Classification, VQA, stitching, etc
- Does not retain state. E.g. you cannot store images.

Previous attempts: NVidia DIGITS

- "DIGITS (the Deep Learning GPU Training System) is a webapp for training deep learning models. "
- Load/create datasets, train models, deploy models.
- Aimed at researchers
- Written in Python/Flask with Torch & Caffe supported

Previous attempts: Visipedia



Taken from Vision of a Visipedia, Perona et. al.

Previous attempts: Visipedia

- Collaborative creation of visual data
- Pre-defined set of concepts E.g. Birds, Trees
- Different type of participants
 - Experts, Annotators, Citizen Scientists, Users, Computer scientists
- Retains state

Previous attempts: VMX.ai

- Underfunded Kickstarter project Circa Jan 2014
- by Tomasz Malisiewicz
- Pre Tensor Flow, Pre Deep Learning
- Allow developers to create real time detectors
- Support for training model

Ongoing attempts

- Scanner by Alex Poms (CMU) & Will Crichton (Stanford)
 - <https://github.com/scanner-research/scanner>
- Kitware Image and Video Exploitation and Retrieval
 - <https://github.com/Kitware/kwiver>
- VISE project by Oxford VGG group
 - <https://gitlab.com/vgg/vise>

Quick recap

- LIRE: limited functionality (Lucene add-on)
- CloudCV: Provides a service, cannot retain “state”
- NVidia Digits: Intended for training not inference
- Visipedia: Intended to be a monolithic deployment

Why now?

- High quality libraries and pre-trained models
 - TensorFlow, PyTorch
 - Inception, SSD, Facenet
 - Flickr LOPQ, Facebook FAISS
- Cheap GPUs (local & cloud)
- Docker enables deployment of complex applications

Relational data : Postgres, MYSQL, SQLite

::

Text, HTML : Lucene/Solr, Elasticsearch

::

Videos & Images : ***Deep Video Analytics***

Provides images & videos,
along with metadata,
annotations

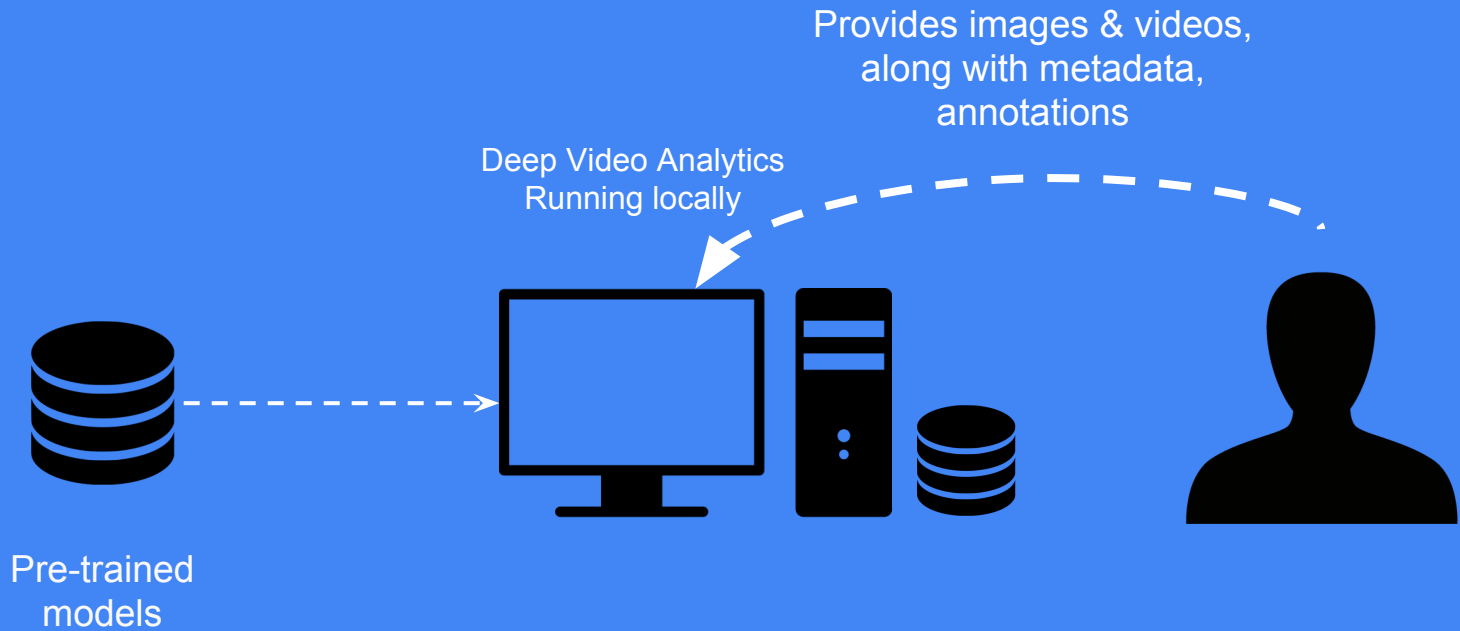
Deep Video Analytics
Running locally

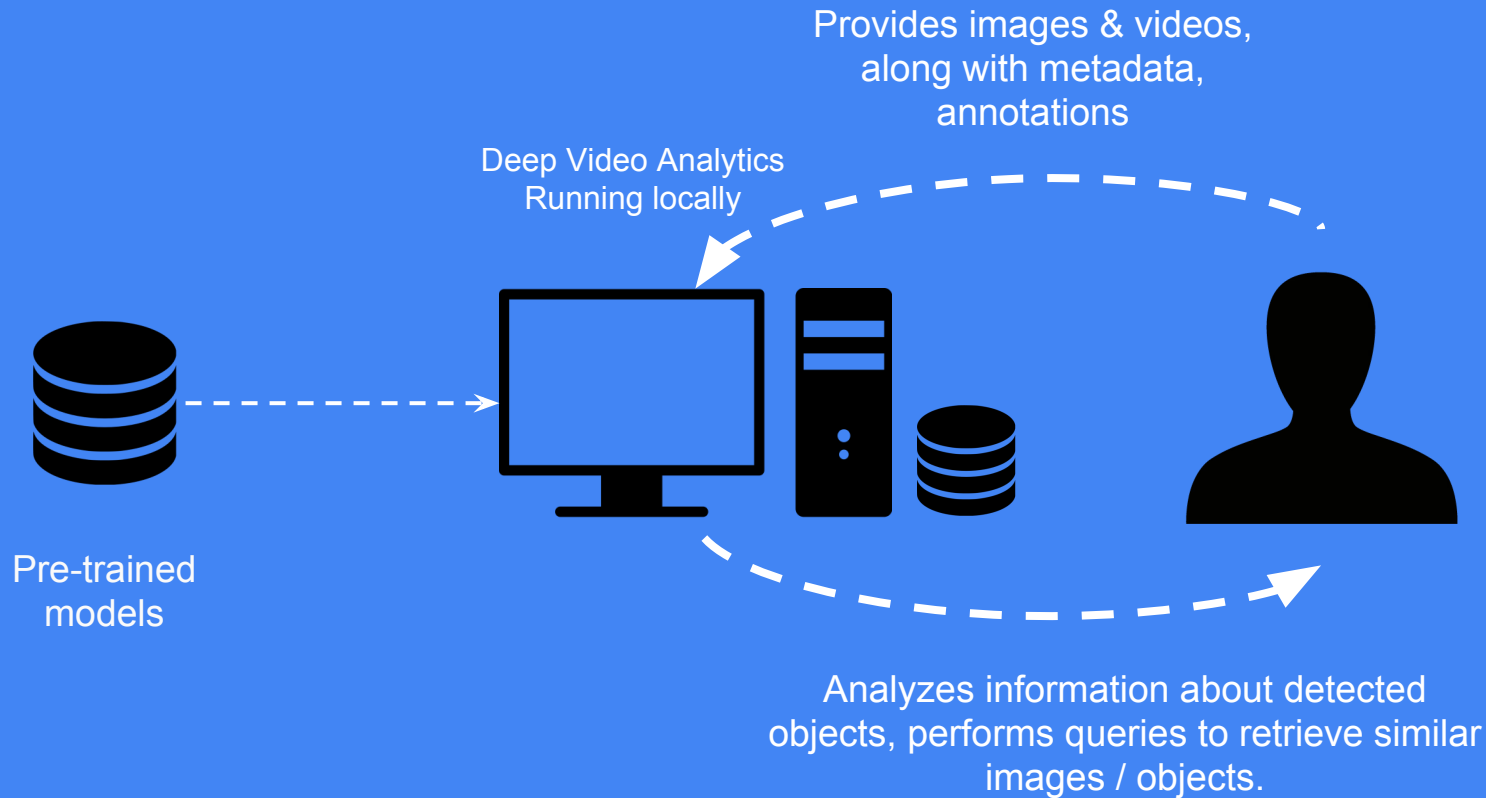


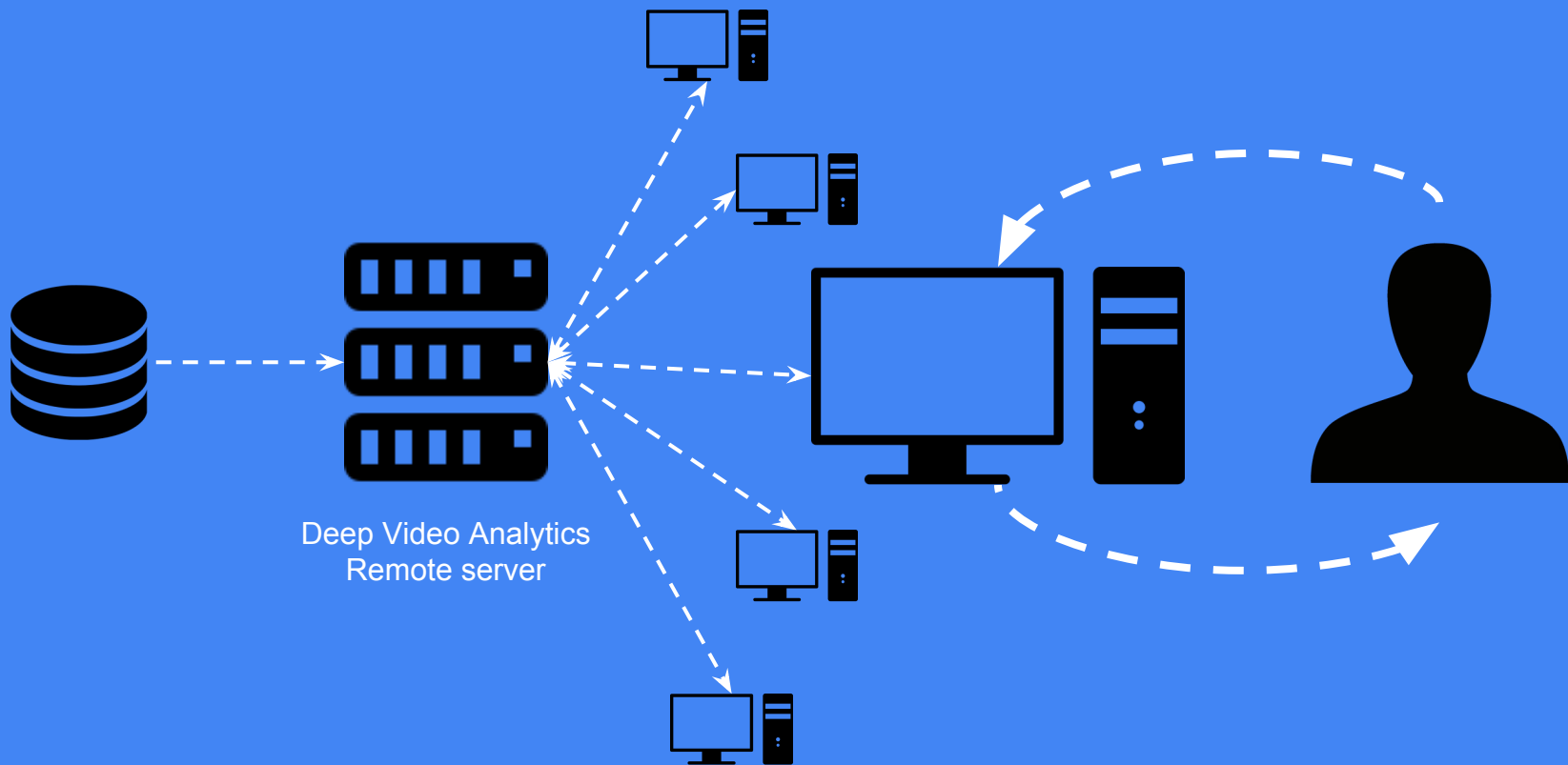
Provides images & videos,
along with metadata,
annotations

Deep Video Analytics
Running locally









Design goals

- Usable by non-researchers
- Visual Search as a “Primary User Interface”
- Users can provide data easily (via upload, youtube-dl, annotation UI etc.)
- Batteries-included approach with an indexing and detection pipeline
 - Tensor Flow Inception v3, VGG-16, Single Shot Detector trained on COCO
 - Face detection / alignment / recognition
 - Deep OCR using CRNN & CTPN. Train new detectors using YOLO+Keras.
- Pre-indexed datasets from different domains can be quickly loaded
- Can be easily customized by developers & researchers.

Technical goals

- Useful without having to write code or config
- Works on machines with and without GPUs
 - Works (albeit slowly) without a GPU, tested on Linode VPS with 8Gb RAM & 4 Cores
- Handles uploads and continuous index updates
- Data can be easily imported, exported and shared
- Can be easily modified by technical users
 - E.g. Adding more operations to processing pipeline
- Can be scaled out by adding more GPUs / Machines

Frameworks & libraries used

- Django, Postgres, Celery, RabbitMQ, Docker, NVidia-Docker
- Tensorflow (primary), PyTorch, OpenCV, FFmpeg, LOPQ & Caffe



What are the core primitives for
Visual Data Analytics?

Visual Data

=

{ Images, Videos, Annotations, Features }

Data & Processing

Data

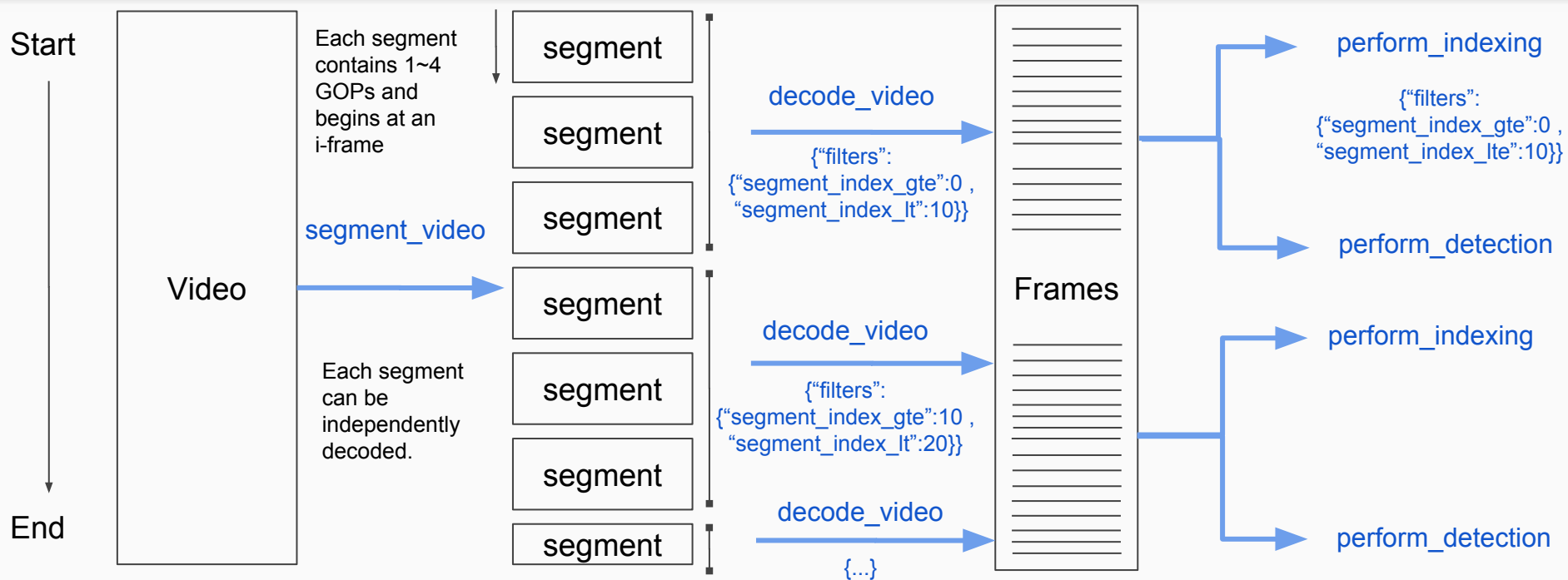
- Video / Segment
- Dataset
- Frame / Image
- Regions over an image
- Tubes over sequence of images
- Feature vectors
- Audio

Processing

- Video Segmentation + Decode
- Image processing
 - Indexing / Detection / Segmentation / Analysis
- Vector processing
 - Retrieve nearest neighbor / Build K-NN graph
- Image transformation
 - Crop / Resize / Align / Apply segmentation mask

Video processing

Parallelized segment + decode pipeline



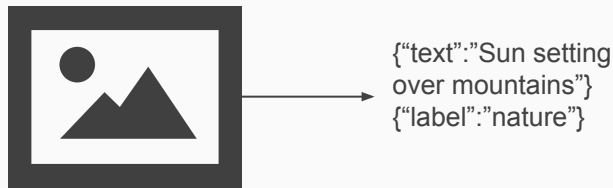
Frame/Region processing operations

Indexing



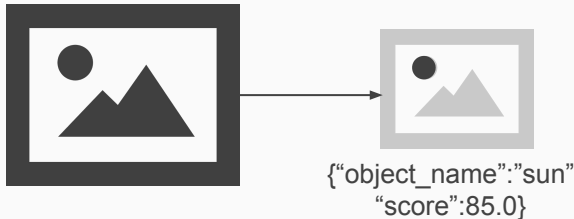
Compute feature vector such as Inception pool, embedding, RGB histogram etc.

Analysis



Analyze image/region and generate metadata (E.g. text description) and/or label

Detection



Detect objects and return bounding boxes

Segmentation



Compute pixel-wise mask using semantic segmentation, superpixels etc.

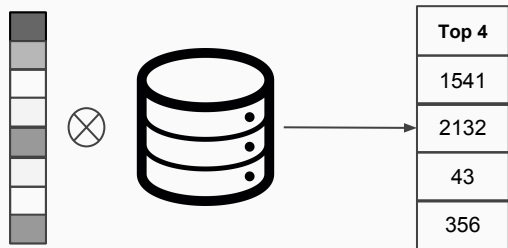
Data & Processing

Key insights

- Different operations have different requirements
 - In terms of number of operations and memory
 - Segmentation > Detection > Indexing / Analysis
- Also different I/O access patterns
 - Detection & Analysis does not requires writing to file system only DB and read
 - Indexing requires writing to filesystem to store computed vectors
 - Segmentation requires writing to filesystem to store computed masks as .png files
- By separating operations we can reason about hardware requirements

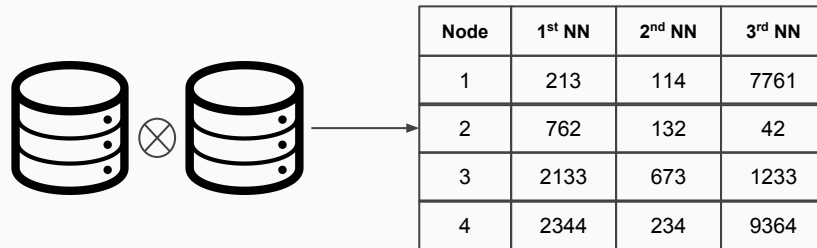
Vector processing operations

Retrieval



Given feature vector find K-Nearest Neighbors

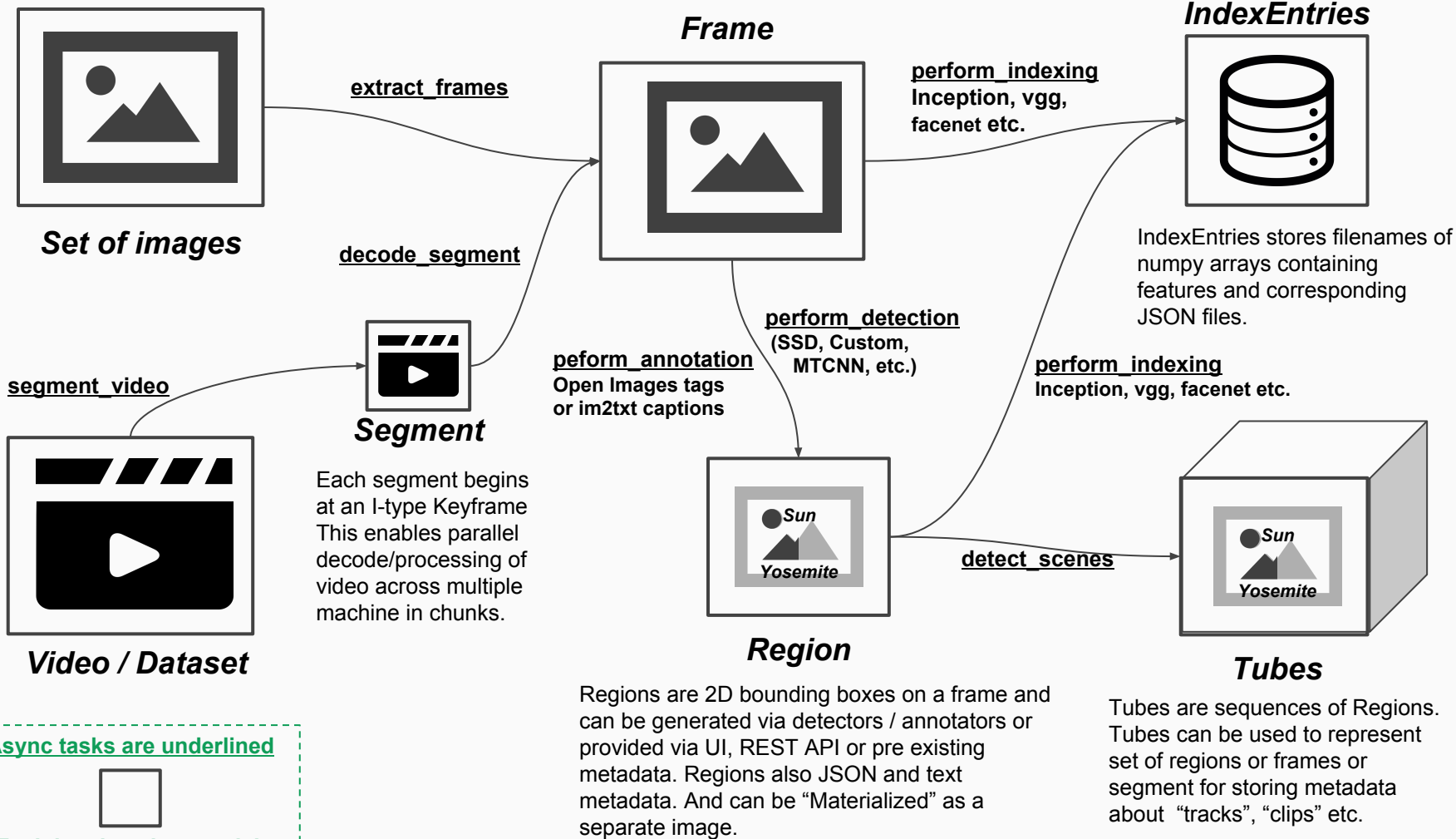
Matching



Given a set of vectors generate K-NN graph

Leverage latest open source implementations for approximate & exact Nearest Neighbors

- Yahoo Locally Optimized Product Quantization (Apache)
- Facebook AI Similarity Search (BSD + **PATENTS restrictions**)



DVAPQL

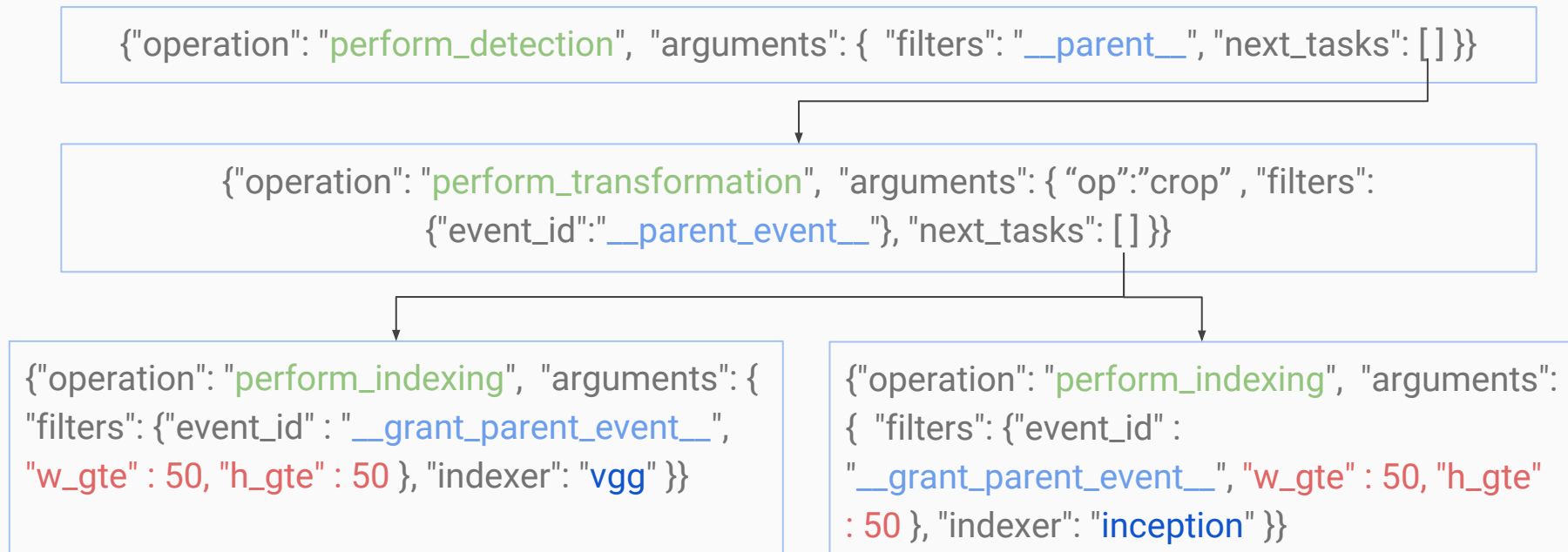
Deep Video Analytics Processing & Query Language

- Specified in JSON
- Launch multiple hierarchical tasks
- Three types of processes
 - Query
 - Retrieve similar images, etc.
 - Process
 - Import video, index images, detect, etc.
 - Schedule
 - Monitor video stream, etc.
- REST API for viewing state & submitting DVAPQL

Example

```
{ "process_type": "V", "tasks": [  
  { "operation": "perform_indexing", ... }  
]  
  
{ "process_type": "Q", "b64_image_data": ".....",  
  "queries": [ { "indexer_query": "perform_indexing", ... }  
]  
  
{ "process_type": "S", "tasks": [  
  { "operation": "ingest_video", ... }  
]}
```

A task based hierarchical processing model



All above tasks run on a specific video / dataset which is not shown for brevity.

Queues for optimal task processing

- Different tasks have different requirements
 - Retrieval / Nearest neighbors: High Memory for storing Index / Approximate index
 - Indexing : GPU for computing embeddings
 - Detection / Segmentation : GPU with higher memory
 - Video decode: GPU optional
 - Crop / Transform / Extract : CPU
- Primitives for Queue management
 - launching queues
 - Monitoring GPU Memory utilization / allocation

Deep Video Analytics

Code organization: dvaapp & dvalib

dvaapp: a django app/project

- Handles UI and data processing
- Data model & Filesystem handling
 - Video, Frame, Detection
 - Query, QueryResult
 - Event, etc.
- Data processing framework using Celery
 - Perform tasks
 - Manage queues
 - Monitor resource use
- Uses dvalib to carry out tasks

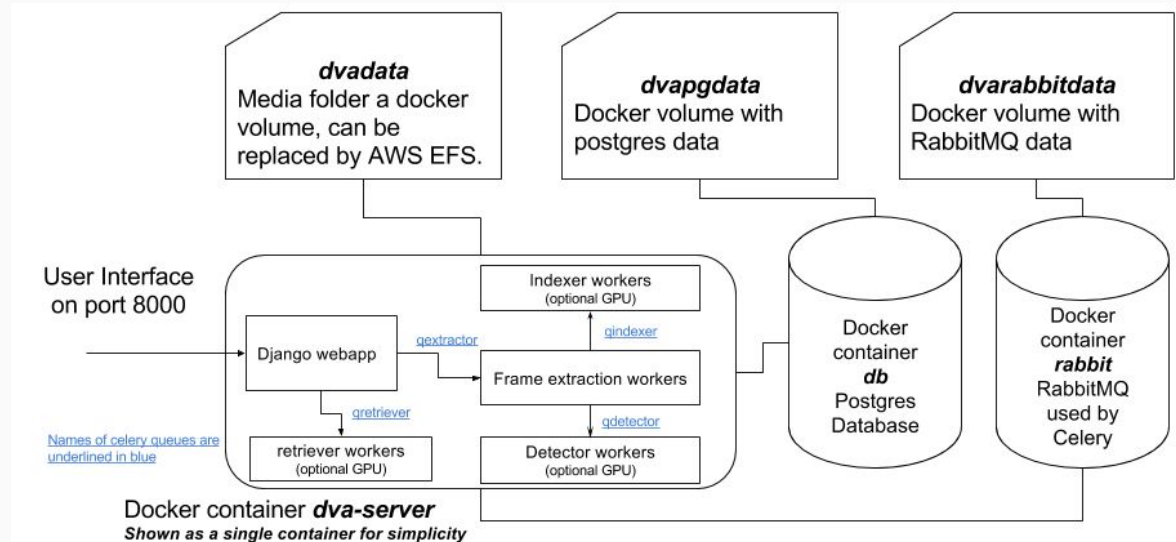
dvalib: library for handling algorithms

- A database & celery agnostic library
- Interface with Tensor Flow & Pytorch for
 - Detection
 - Indexing
 - Segmentation

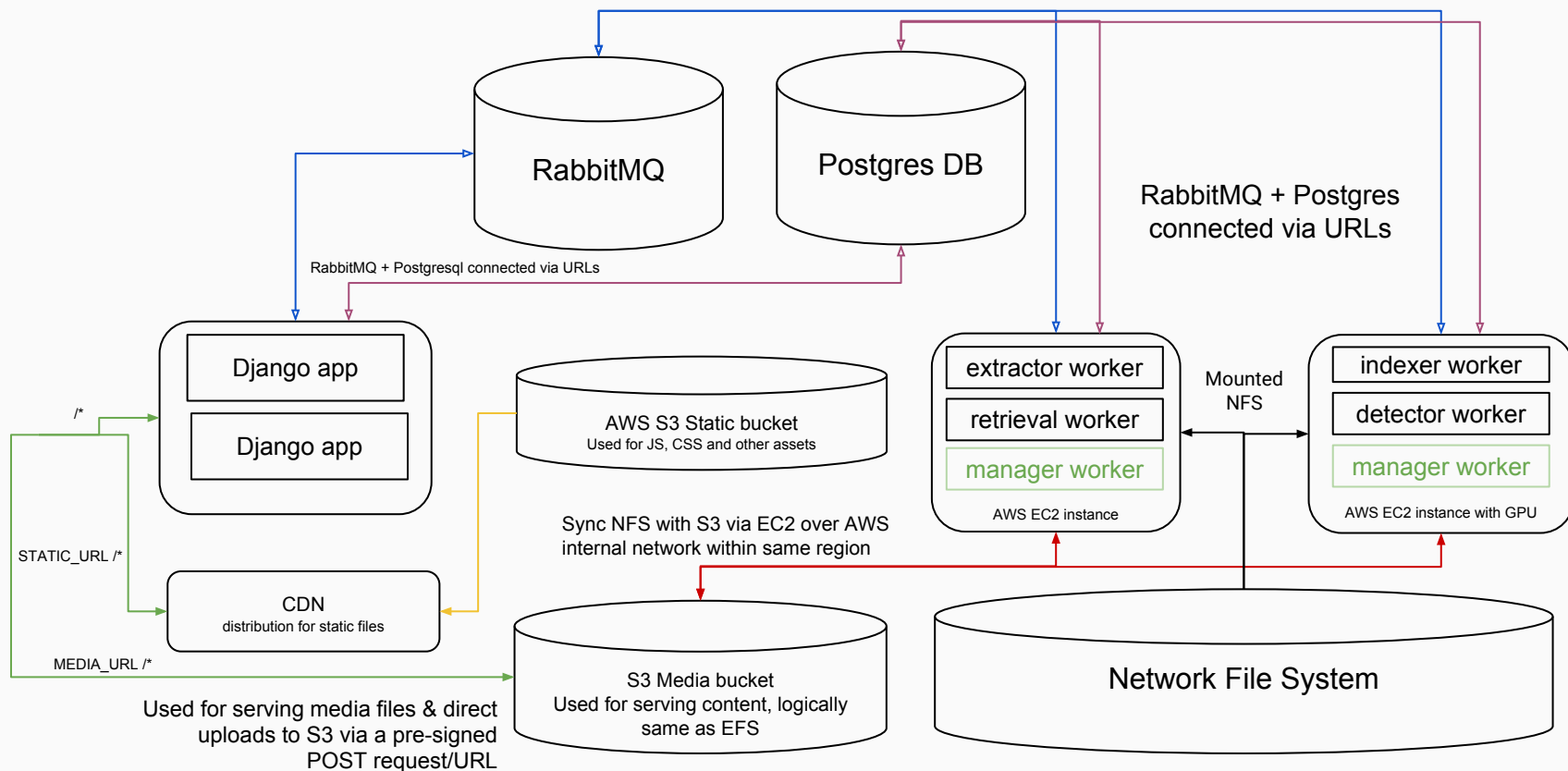
Emulating datacenter on a machine

Docker, Docker-compose & Nvidia-docker

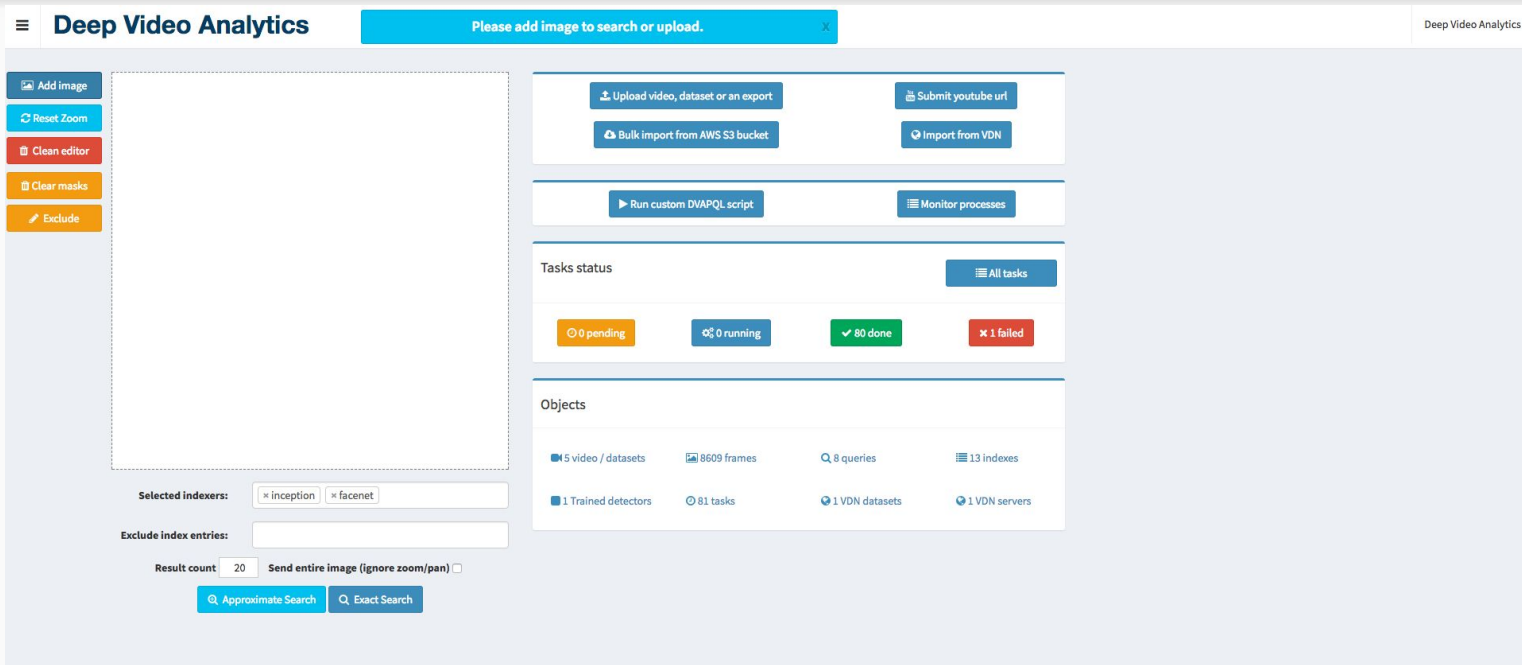
Docker enables same codebase across all configurations {a laptop, multi-GPU machine, datacenter} .



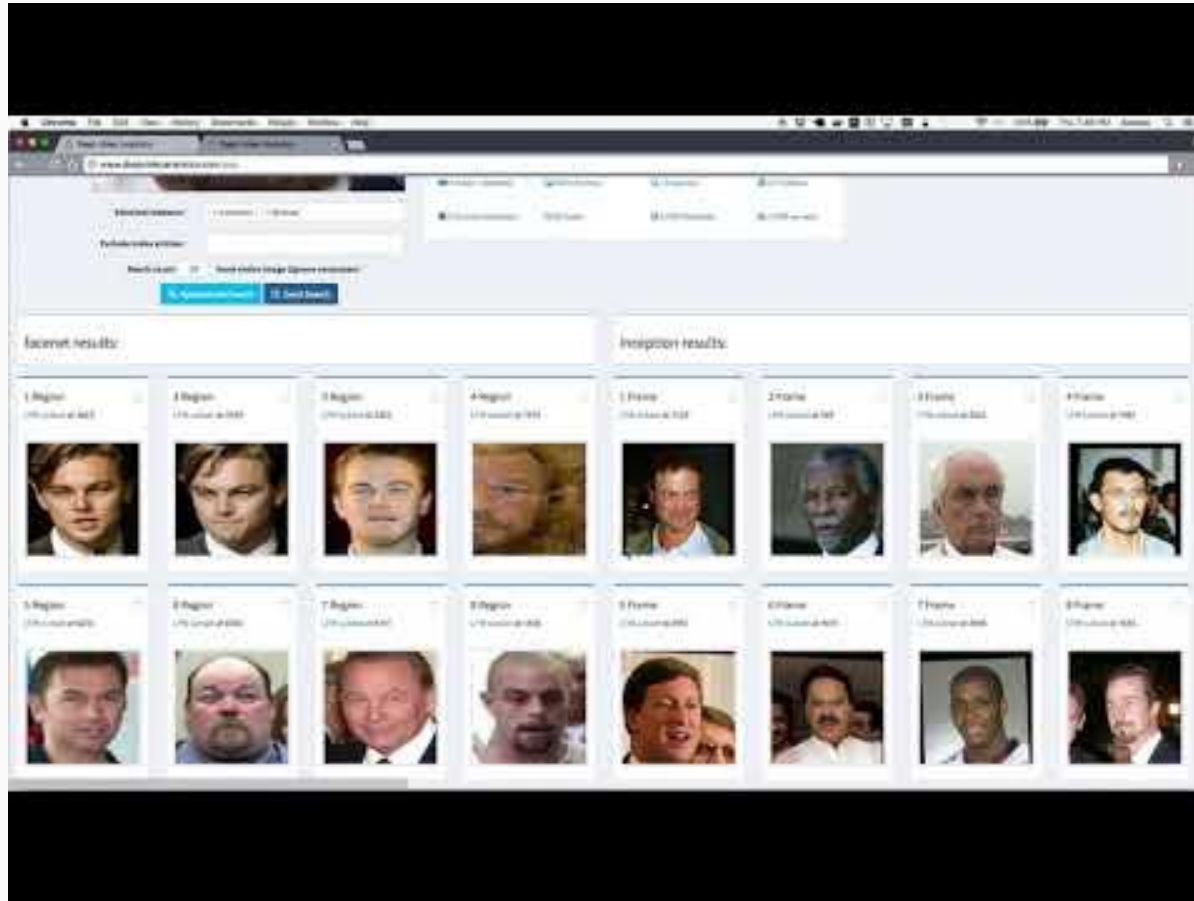
Scalability with distributed architecture



User Interface



Latest version beta, 17th August 2017



7th April 2017

Deep Video Analytics

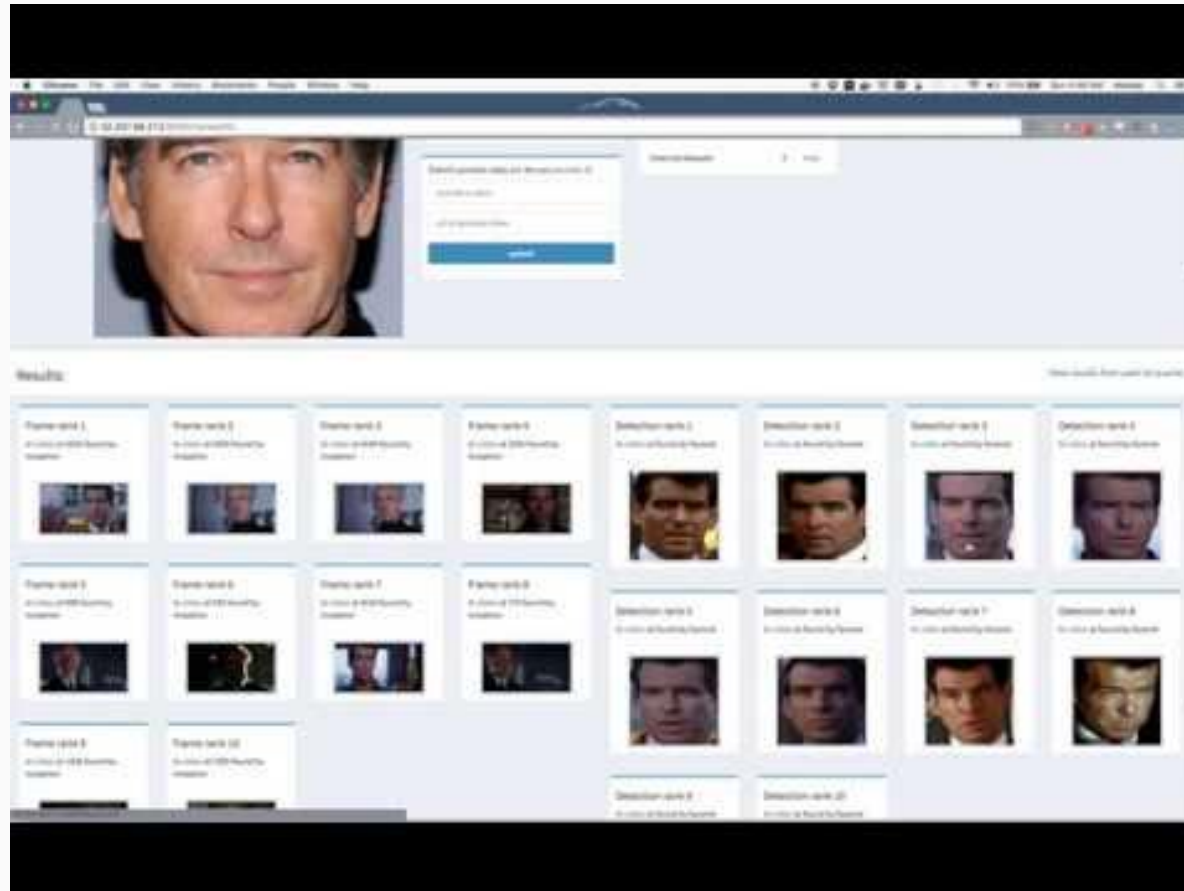
134.47 (seconds) of ingested content, January, 2021 No live feed 2022 1440x1024, 30FPS, 119.1MB/s

ID	Class	Confidence	x1	y1	x2	y2
111	Class (string)	0.9944 (float)	44	15	461	165
1000000000	Class (float)	0.9944 (float)	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165
1000000000	Class	0.9944	44	15	461	165

Detected objects:

1000000000

15th March 2017



People : Facebook

::

Code : Git / GitHub, GitLab

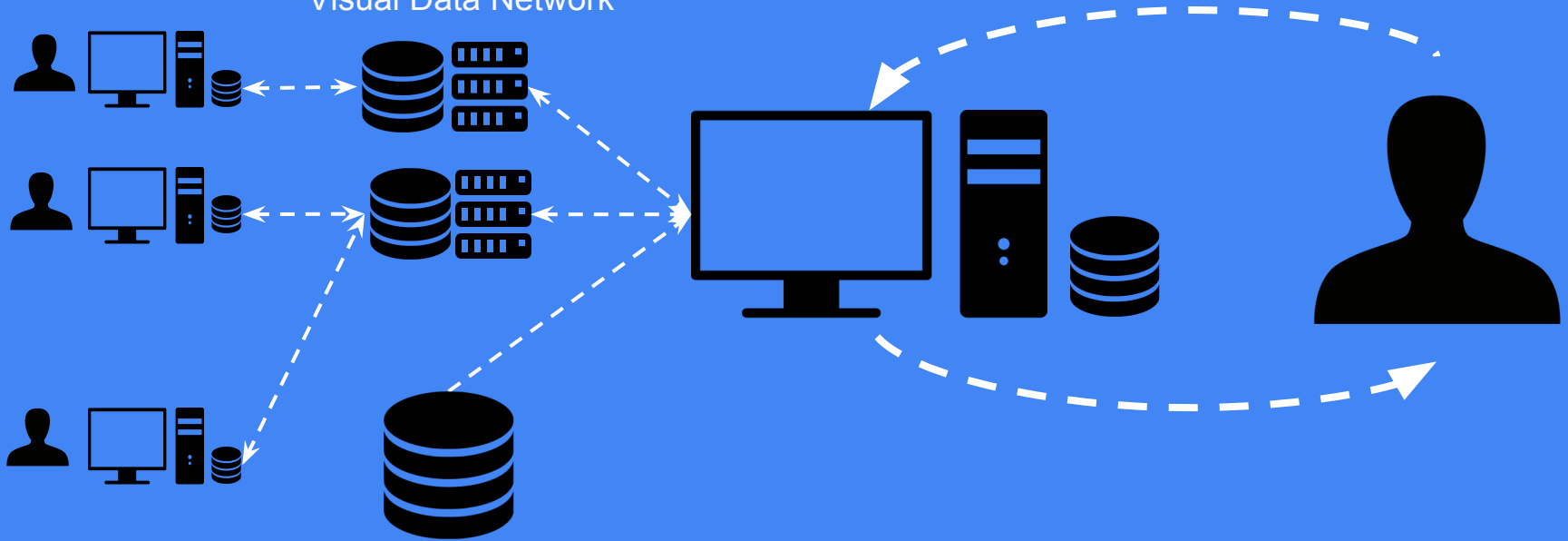
::

Visual Data: ***Visual Data Network***

Sharing data using Visual Data Network

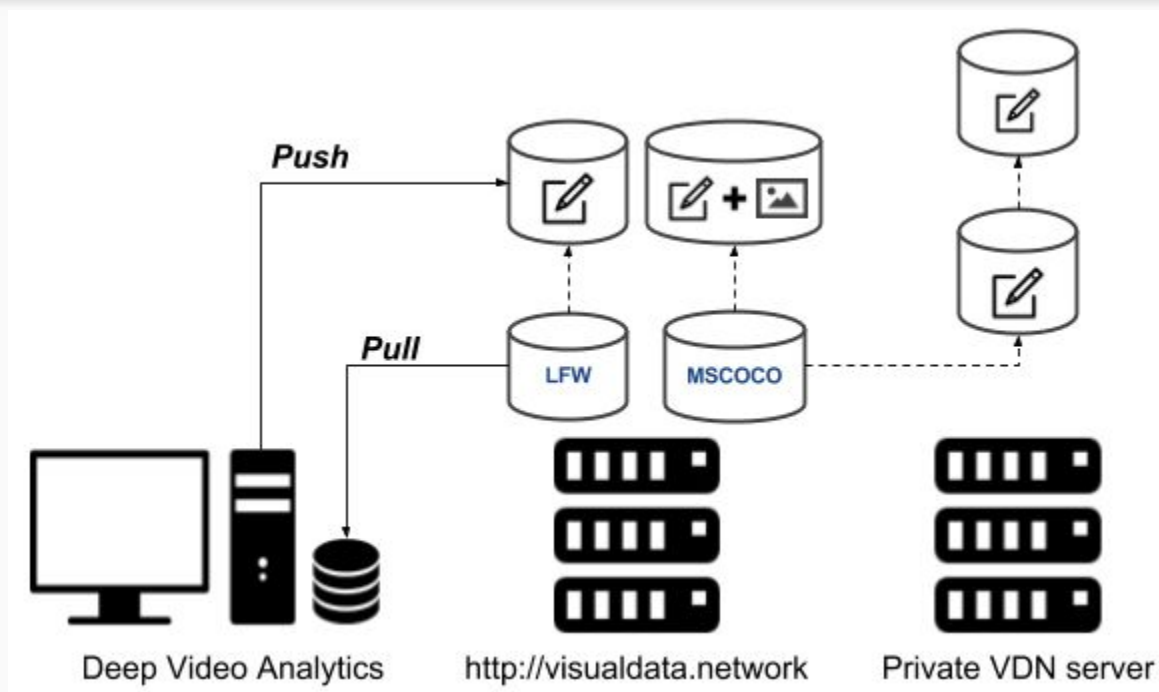
Import & export new datasets / annotations
share with other users

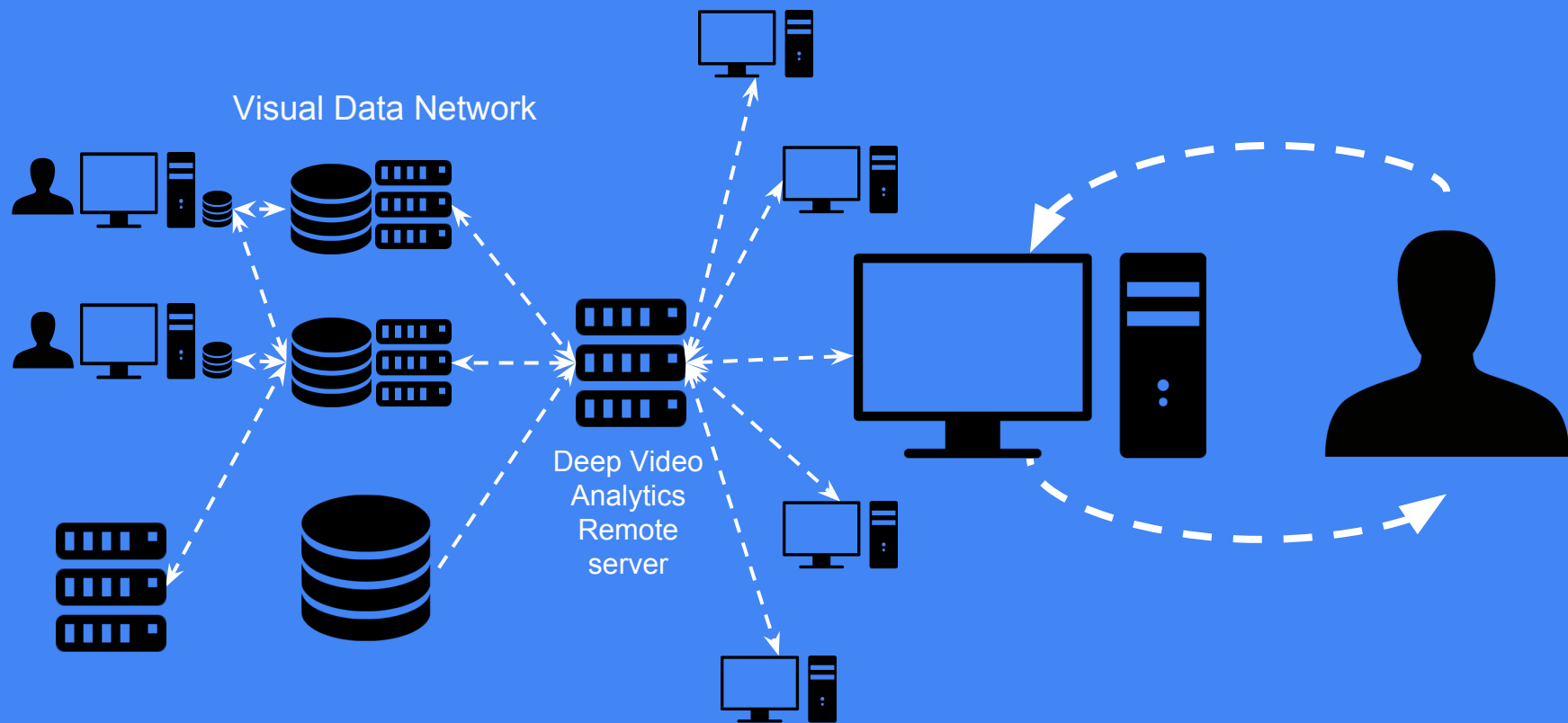
Visual Data Network



Visual Data Network enables seamless sharing

Push, Pull video / dataset, Annotations, just like you would with GitHub





Open questions:

A work in progress

- How to effectively manage GPU memory & utilization?
- How to balance fast/static vs slow/dynamic indexes?
- How to learn continuously from annotations/feedback?
- How to minimize storage requirements via compaction?
- How to enable Real time processing?

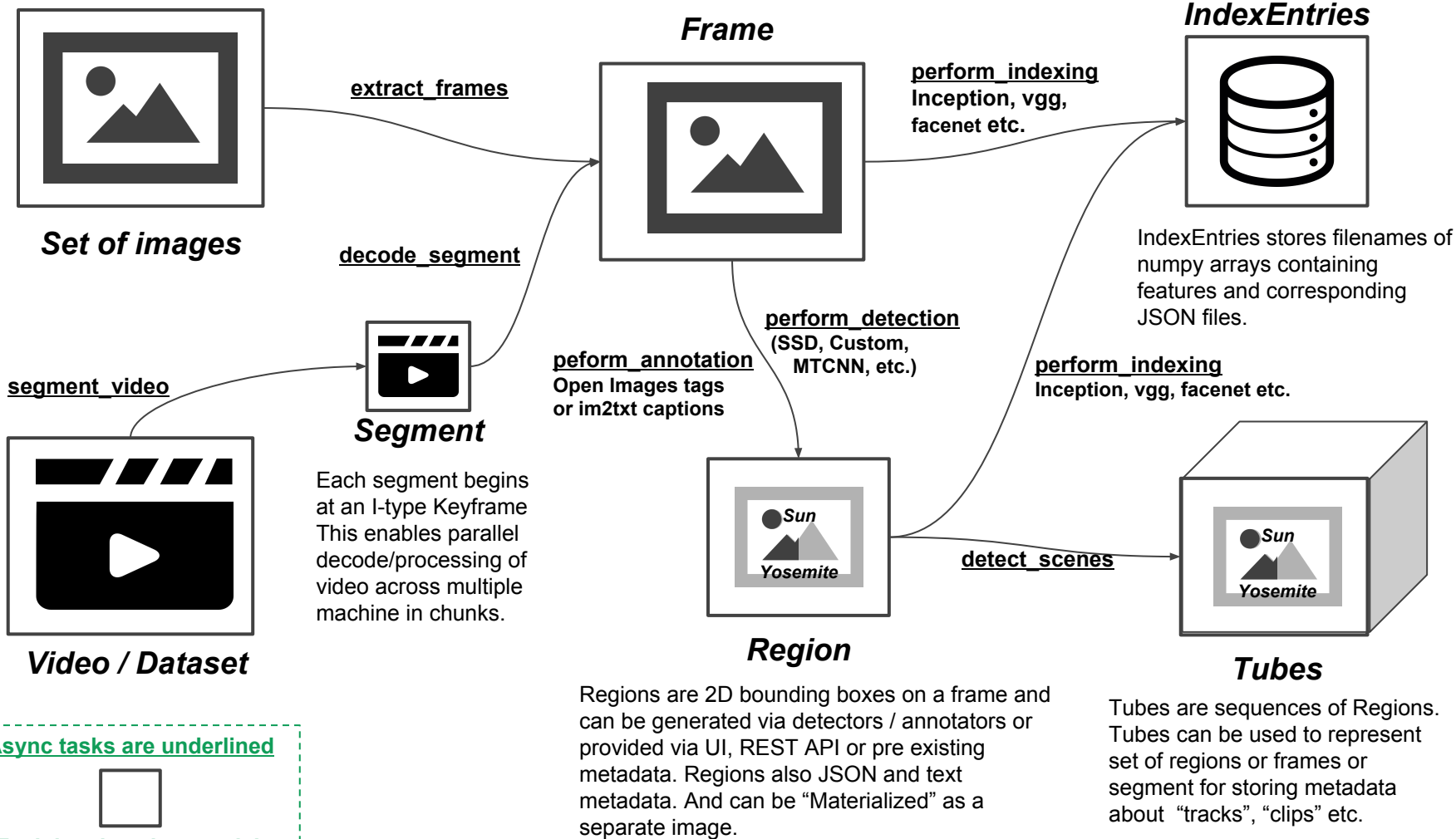
Thanks!

Contact me:

akshayubhat@gmail.com

www.akshaybhat.com





Distributed processing using hierarchical tasks

