# Visual Intelligence Platform

## Deep Video Analytics + Visual Data Network

Akshay Bhat
Cornell Tech, Cornell University.

# An overview of computer vision research by Tomasz Malisiewicz

# Quick summary

Sift, Graph Cuts

↓

HOG, DPM

↓

Deep Learning

↓

?

Caltech 101, Matlab, OpenCV

↓

VOC, Imagenet, Caffe, Theano

↓

?

# Numerous high quality libraries

- OpenCV
- ROS
- Caffe
- Theano
- Torch

- Tensor Flow
- CNTK
- MXNET
- Torch
- deeplearn.js

# 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

# A deluge of datasets!

- VideoNet
- Yahoo Flickr Creative Commons 100M
- ViCom
- Visual Genome
- YouTube-BoundingBoxes
- Youtube 8M

- imSitu by AllenAI
- Charades by Allen AI
- Udacity car dataset
- KITTI
- Caltech, INRIA, ETH Pedestrians
- Stanford Drone Dataset
- COCO text

We are reaching a stage where

Number of datasets ≅ Number of research groups

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

What is hidden in plain sight?

We need a platform which seamlessly combines

Data + Models + User Interface

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

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

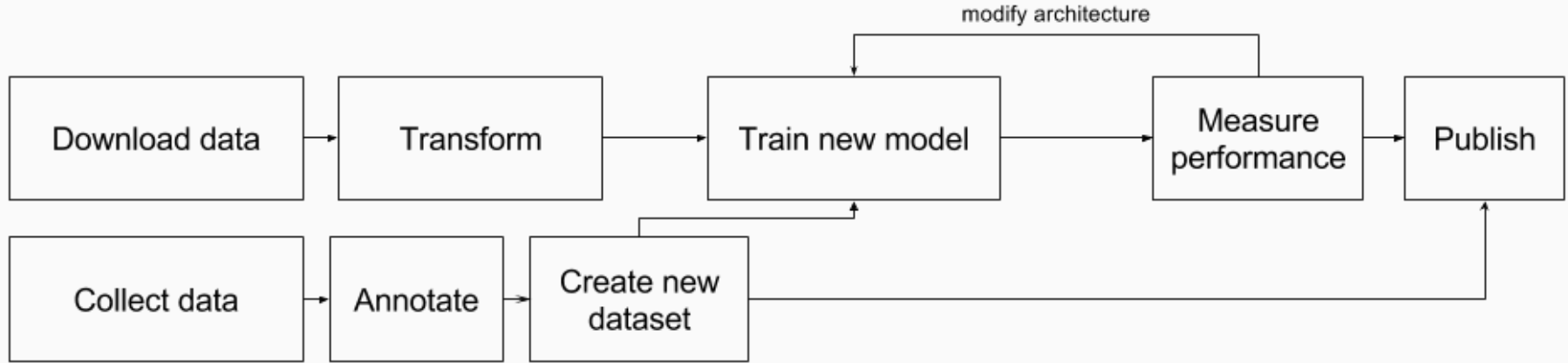Visual Data
=
{ Images, Videos, Annotations, Features}

Relational data : Postgres, MYSQL, SQLite
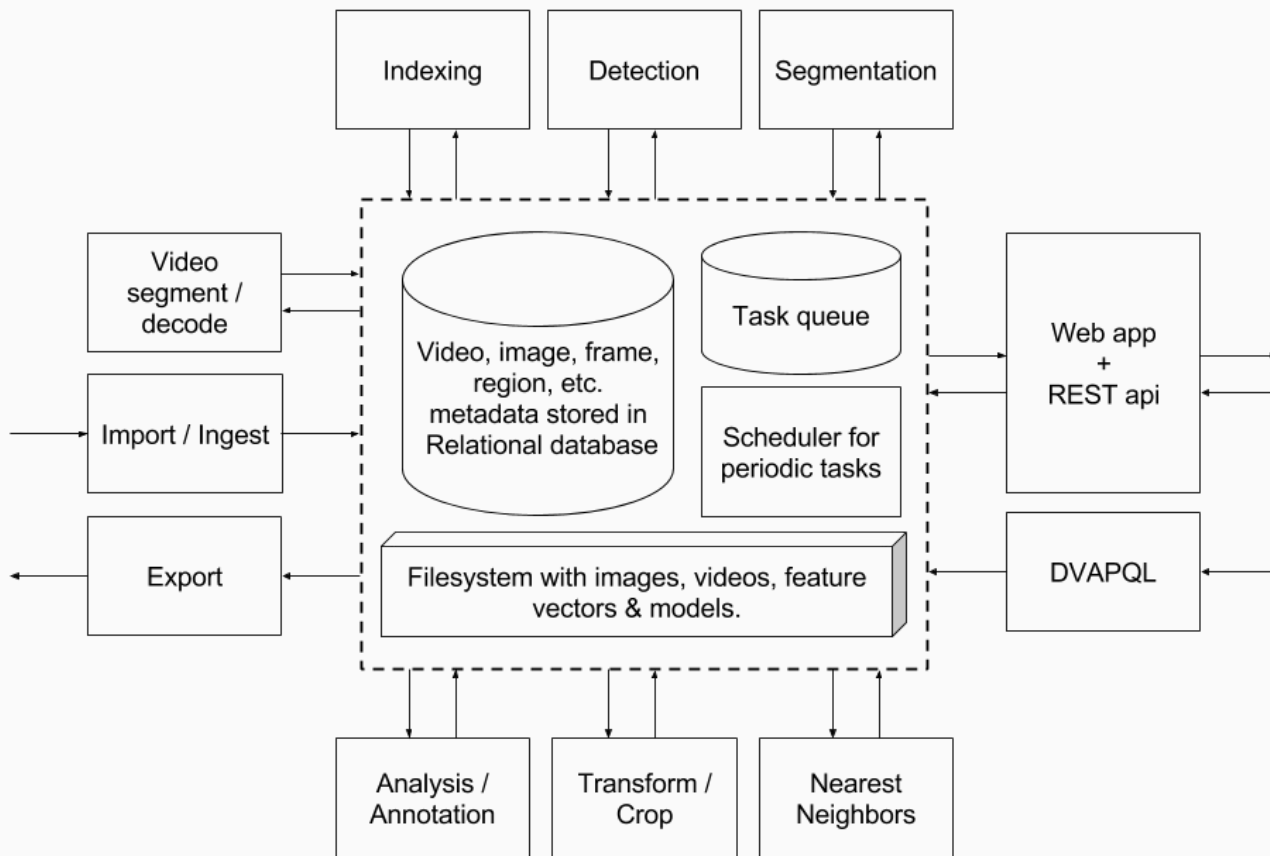::
Text, HTML : Lucene/Solr, Elasticsearch
::
Videos & Images :  _____

# Model-centric

# Model-centric to **Data-centric**

# 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.
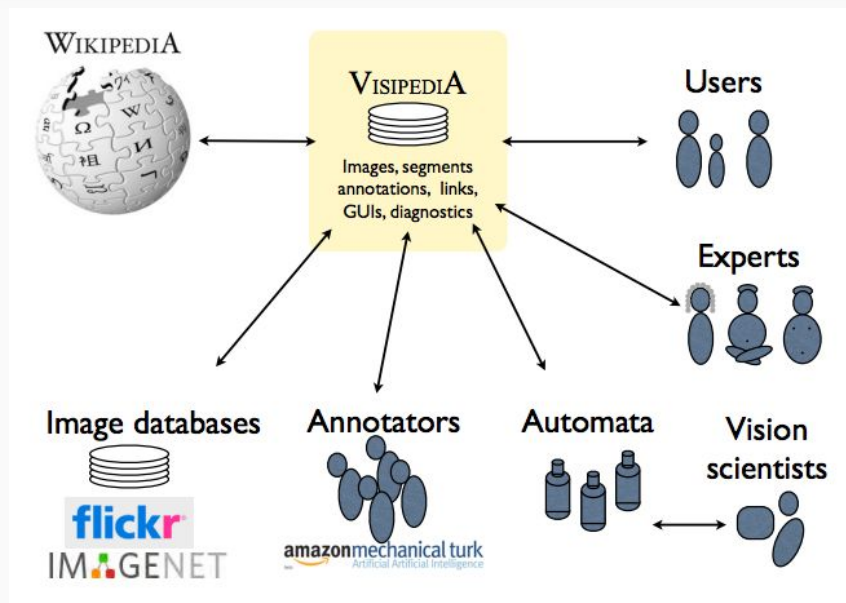
# Previous attempts: CloudCV

- Large Scale Distributed Computer Vision as a Cloud Service

- Support for OpenCV, Graphlab, Cafe

- Image Classification, VQA, stitching, etc

- Does not retains 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

# 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

# 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

# Why now?

- High quality libraries and pre-trained models
  - TensorFlow
  - 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 :  _____

Relational data : Postgres, MYSQL, SQLite

::

Text, HTML : Lucene/Solr, Elasticsearch

::

Videos & Images :  *Deep Video Analytics*

People : Facebook, MySpace
::
Code : Git / GitHub, GitLab
::
Visual Data: *Visual Data Network*

Relational data : SQL

::

Text, HTML : inverted word index, Page Rank

::

Videos & Images : *Approximate Nearest Neighbor*

Provides images & videos,
along with metadata,
annotations

Deep Video Analytics
Running locally

Provides images & videos,
along with metadata,
annotations

Deep Video Analytics
Running locally

Provides images & videos,
along with metadata,
annotations

Deep Video Analytics
Running locally

Pre-trained
models

Provides images & videos,
along with metadata,
annotations

Deep Video Analytics
Running locally

Analyzes information about detected
objects, performs queries to retrieve similar
images / objects.

Pre-trained
models

# 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

# Flexible deployment: local & remote server



Visual Data Network

Deep Video
Analytics
Remote
server

# 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 & technologies used

- Django, Postgres, Celery, RabbitMQ, FFmpeg, Docker

- Tensorflow (primary), Torch, OpenCV & Caffe

What are the core primitives for Visual Data Analytics?

# Data & Processing

## Data

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

## Processing

- Video Segmentation + Decode
- Indexing
  - Compute features for a region / image
- Detection
  - Detect objects in an image or a region
- Annotation / Analysis
  - Generate a label/metadata given a video, image, region, segment or a tube.
- Transformation
  - Generate a new image/region or tube from existing one. (e.g. segmented object stored as .png file)

**Frame**

**IndexEntries**

**extract_frames**

**perform_indexing**
Inception, vgg, facenet etc.

IndexEntries stores filenames of numpy arrays containing features and corresponding JSON files.

**decode_segment**

Set of images

**perform_detection**
(SSD, Custom, MTCNN, etc.)

**peform_annotation**
Open Images tags or im2txt captions

**perform_indexing**
Inception, vgg, facenet etc.

**segment_video**

**Segment**

Each segment begins at an I-type Keyframe This enables parallel decode/processing of video across multiple machine in chunks.

**Sun**
**Yosemite**

**detect_scenes**

**Sun**
**Yosemite**

**Video / Dataset**

**Region**

Regions are 2D bounding boxes on a frame and can be generated via detectors / annotators or provided via UI, REST API or pre existing metadata. Regions also JSON and text metadata. And can be "Materialized" as a separate image.

**Tubes**

Tubes are sequences of Regions. Tubes can be used to represent set of regions or frames or segment for storing metadata about "tracks", "clips" etc.

Async tasks are underlined

*Each box is a data model*

# DVAPQL
## Deep Video Analytics Process & Query Language

- Specified as JSON
- Three type of scripts
  - Process
  - Query
  - Ingest
- Launch multiple tasks
- Monitor & Wait on tasks
- Use REST API for viewing state & annotation
- Use DVAPQL for launching tasks

Example

{ "process_type" : "V", "tasks": [

{"operation":"perform_indexing", ... ]}


{ "process_type" : "Q", "b64_image_data":".....",

"queries": [ {"indexer_query":"perform_indexing", ...

]}


{ "process_type" : "I", "tasks": [

{"operation":"ingest_video", ... ]}

# A task based flexible processing model

{"task_name": "perform_detection",  "arguments": {  "filters": "__parent__", "next_tasks": [ ] }}

{"task_name": "crop_regions",  "arguments": {  "filters": {"event_id":"__parent_event__"}, "next_tasks": [ ] }}

{"task_name": "perform_indexing", "arguments": {  "filters": {"event_id" : "__grant_parent_event__", "w_gte" : 50, "h_gte" : 50 }, "indexer": "vgg" }}

{"task_name": "perform_indexing", "arguments": {  "filters": {"event_id" : "__grant_parent_event__", "w_gte" : 50, "h_gte" : 50 }, "indexer": "inception" }}

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

# Parallelized video processing segment + decode pipeline

# Emulating datacenter on a machine
## *Docker, Docker-compose, Nvidia-docker*

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

# 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
  - Extract frames / process video
  - Perform indexing
  - Perform detection
- Uses dvalib to carry out tasks

**dvalib:** library for handling algorithms

- A database & celery agnostic library

- Interface with Tensor Flow & Pytorch for
  - extraction
  - detection
  - indexing

# User Interface:
## Search across frames + detections (faces, etc.)

# Demo Version Alpha 1, 15th March 2016

# Demo Version Alpha 2, 7th April 2017

# Open questions:
## A work in progress

- How to rank results using auxiliary information?
- How to balance fast/static vs slow/dynamic indexes?
- How to incorporate text data extracted from images?
- Learning from annotations?
- Real time plug-in that bypasses queue based system?
- An Android / iOS frontend app for data acquisition?

# Thanks!

Contact me:

akshayubhat@gmail.com
www.akshaybhat.com