

PyctureStream

Experts for distributed Image Processing

Imaginary Sales Pitch
by Holger Büch & Marcus Fixel

Master „Data Science and Business Analytics“
Hochschule der Medien, Stuttgart
Modul „BI- and Big-Data-Architectures“
Prof. Dr. Peer Küppers

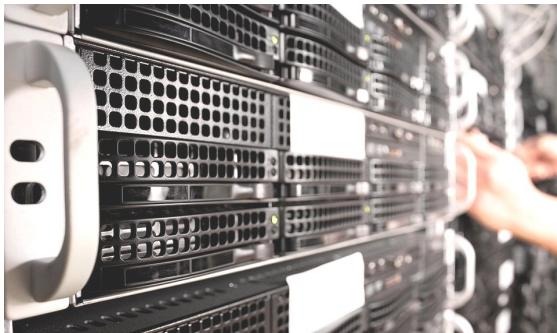
AGENDA

The next 30 Minutes

- 1** Introduction
- 2** Business Case
- 3** Designing the Architecture
- 4** Reference Architecture & Implementation
- 5** Demo of the Prototype
- 6** Wrap Up

INTRODUCTION

PyctureStream - Who are we?



EXPERTS FOR BIG DATA

Two Founders & Consultancy-Team

Academical Research Projects

Experience in building
Real World Applications



SPECIALIZED ON VISUAL DATA

Streaming Images or Videos

Applying advanced Analytics
like Deep Learning

Scalable distributed Architectures



YOUR PARTNER FOR YOUR NEEDS

Project Planning

Requirements Engineering

Architecture &
Implementation

INTRODUCTION

What we can offer?

ESCORTING YOUR BUSINESS

We don't leave the ship, when it's set up. We help you sailing to success.



CUSTOMIZATIONS FOR YOUR NEEDS

We identify your requirements and adept our Solution to your Business.



REFERENCE ARCHITECTURE

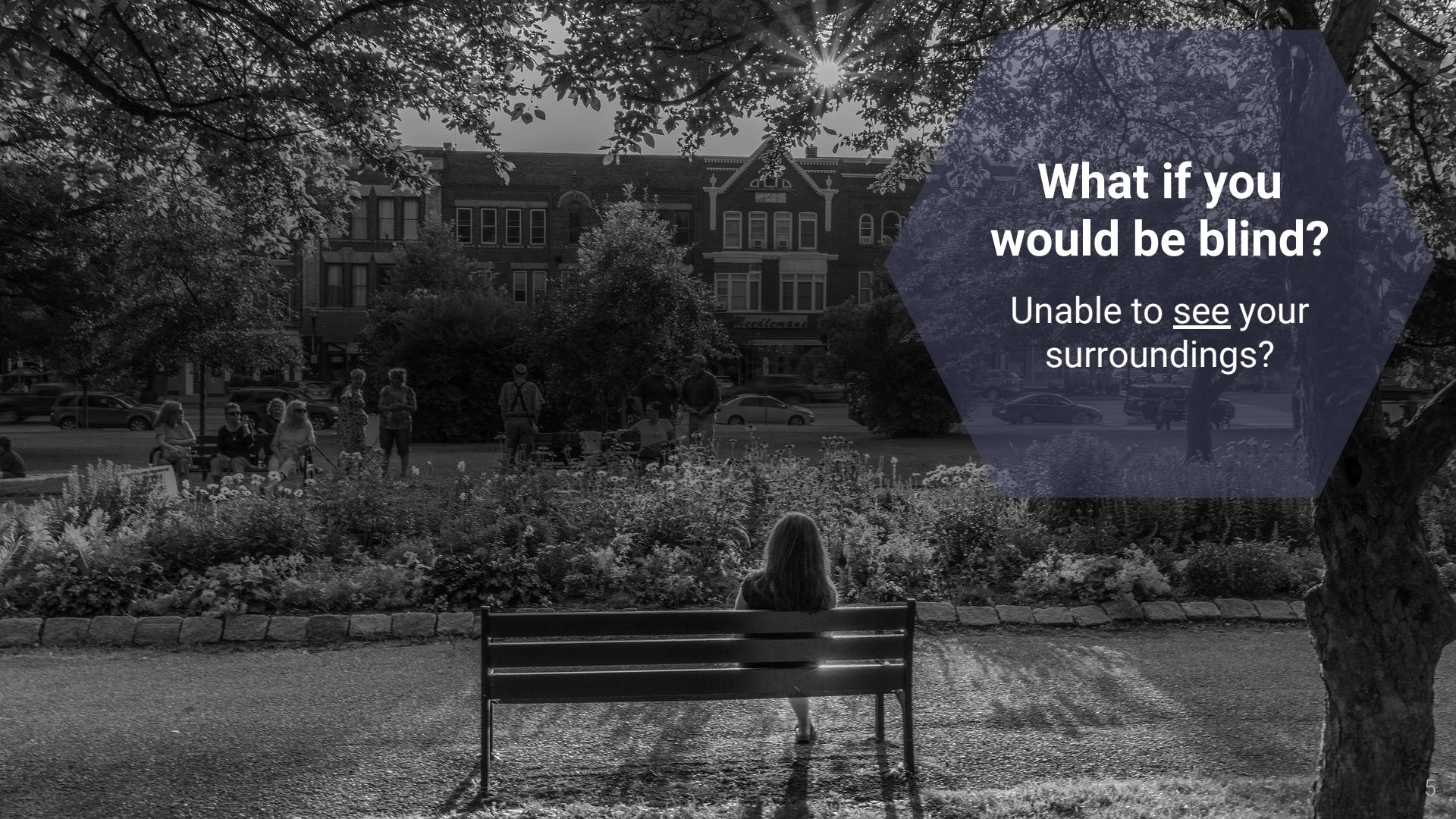
The Blueprint for your Big Data Architecture, proven in multiple Projects.



EXPERIENCED BIG DATA ARCHITECTS & DEVELOPERS

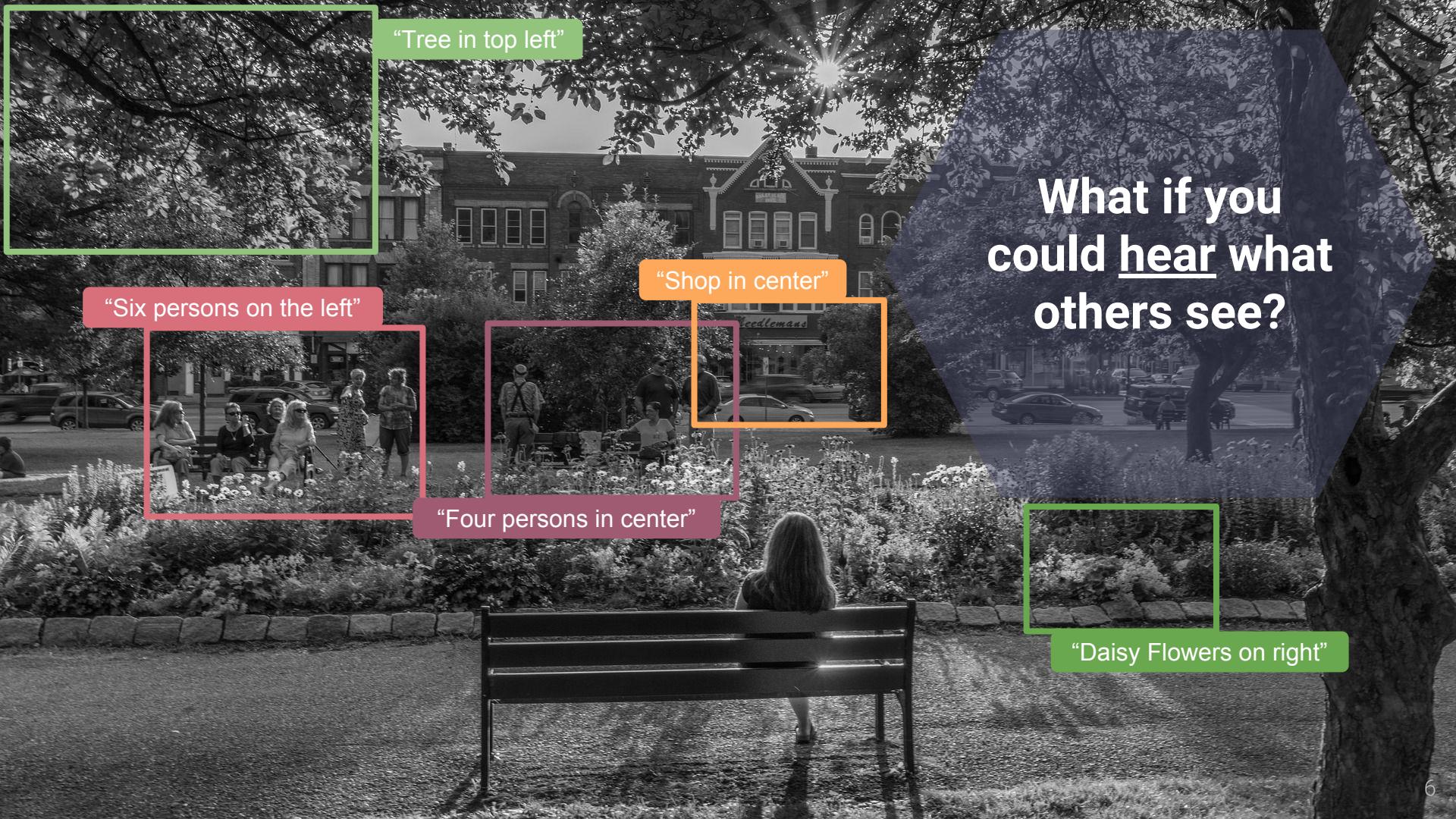
We have Data Scientists and Data Engineers with the right Skill Set.





**What if you
would be blind?**

Unable to see your
surroundings?



What if you
could hear what
others see?

"Tree in top left"

"Six persons on the left"

"Shop in center"

"Four persons in center"

"Daisy Flowers on right"

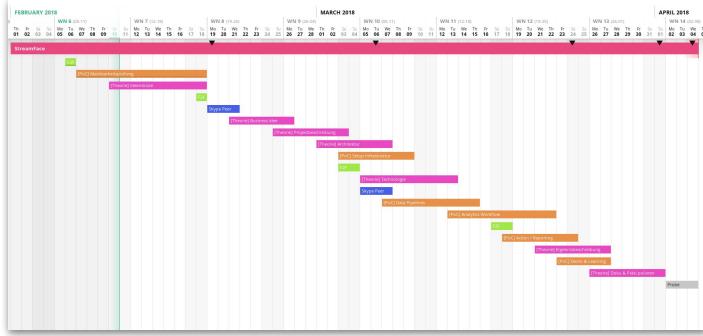
BUSINESS CASE

Business Model Canvas

Key Partners	Key Activities	Value Proposition	Customer Relations	Customer Segment
Nonprofit Organizations	Image Analysis	More Information about Surroundings More Independence Fast & easy to use	Subscriptions Technical Support	Visually impaired People Blind People
	Key Resources	Data Protection	Channels	
	Infrastructure Development Experts		Smartphone App Online Media	
Cost structure		Revenue Stream		
Operation of Data Center Development Costs Sales and Marketing Costs		Donations Subscription Fees Partners		

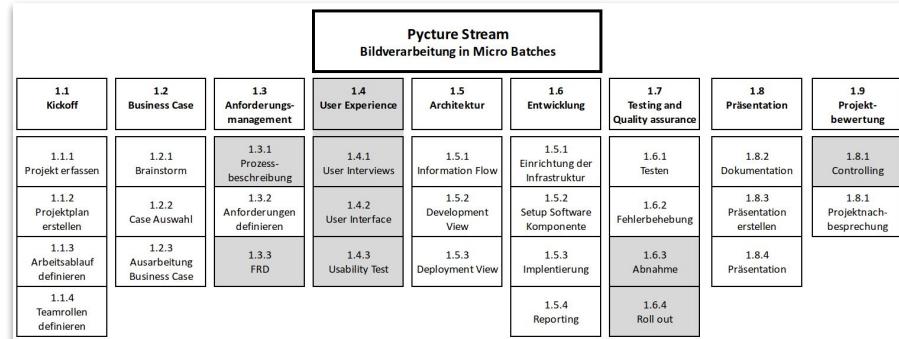
BUSINESS CASE

Project Planning



Schedule-oriented Planning
e.g. using Gantt-Charts

Activity-oriented Planning
e.g. using Work Breakdown Structures



Business Requirements

CAMERA-STREAM PER DAY

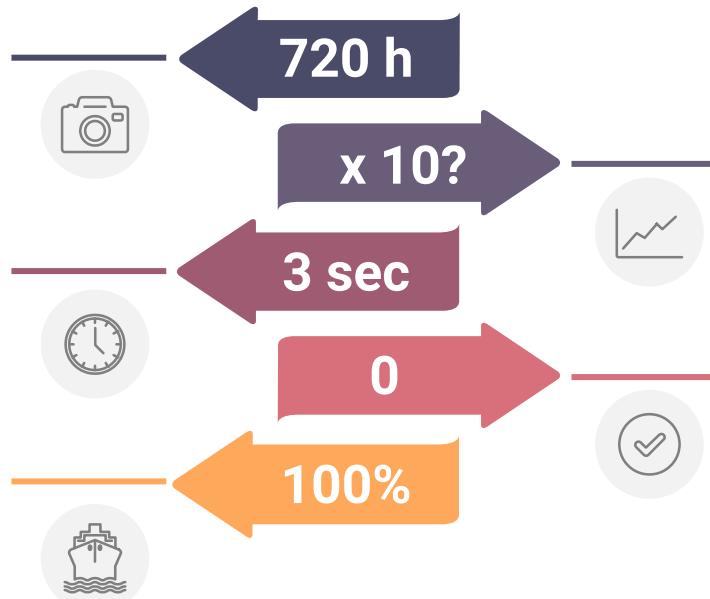
~4000 Customers using the App
for 10 min per Day

INSTANT RESULTS

Timespan from taking the picture
to hearing the results

ENABLEMENT

Enable the Company to drive the
System and avoid Vendor Lock In.



SCALABILITY

Flexible Scaling in the first three Years,
depending on Usage.

PRIVACY CONCERNs

No images or videos of customers
should be used or stored. Ever.

DESIGNING THE ARCHITECTURE

Technical Requirements



Business Requirements
&
Technical Expertise



Reference Architecture

DESIGNING THE ARCHITECTURE

Basic Design Decisions

Big Data Architecture?

- Volume, Velocity, Veracity, Variety

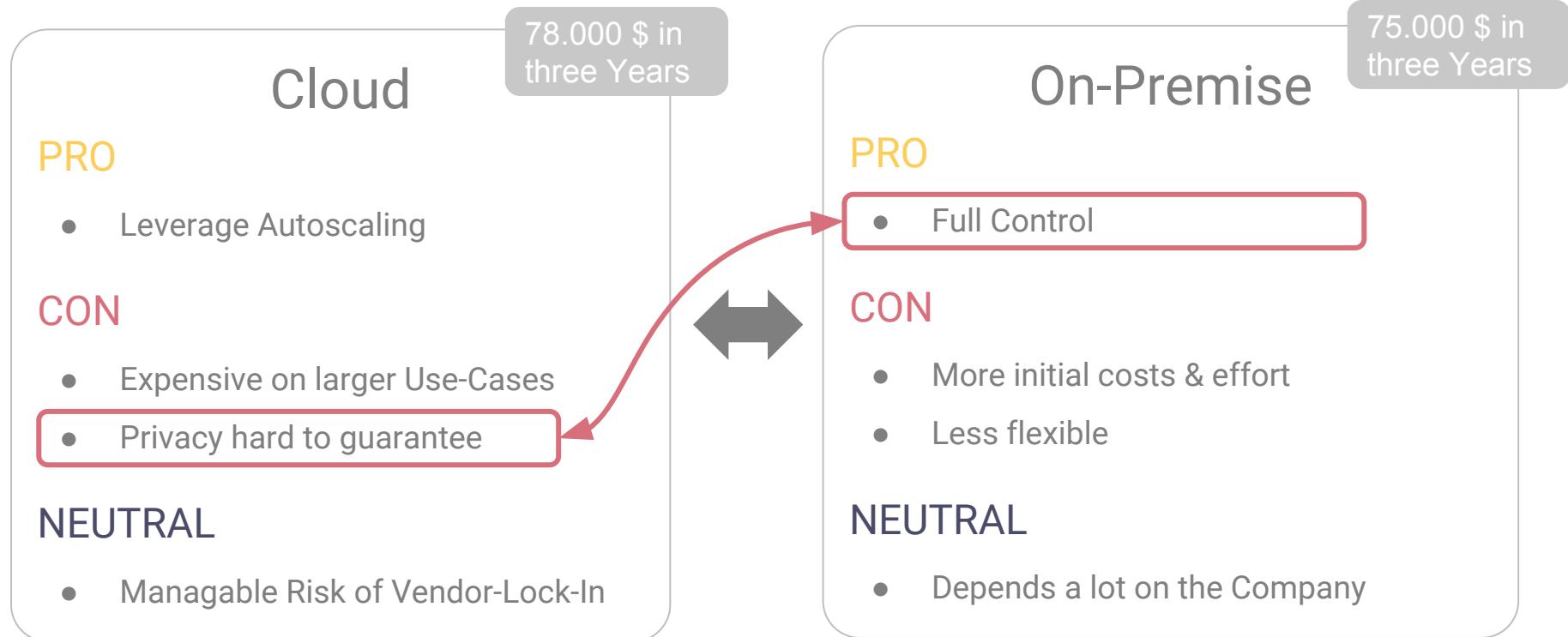


Distributed Architecture?

- Flexible Scaling, Preconditions



Choosing the Platform



Deciding on Architecture Model

Lambda

PRO

- Enabling Re-Use of Images

CON

- More complex & effort
- Batch Layer not needed

NEUTRAL

- Enables Batch & Stream Approaches

Kappa

PRO

- Focused on Streaming
- Single Technology Stack

CON

- Less widely used

NEUTRAL

- Skills for Streaming Analytics required

Selecting the Software

Kafka, Spark & TensorFlow

PRO

- Widely known, used & supported
- Very flexible
- Solid Choice

NEUTRAL

- Lots of Features not needed in all Use-Cases

Possible Alternatives

Streaming

- RabbitMQ, Apache ActiveMQ

Processing

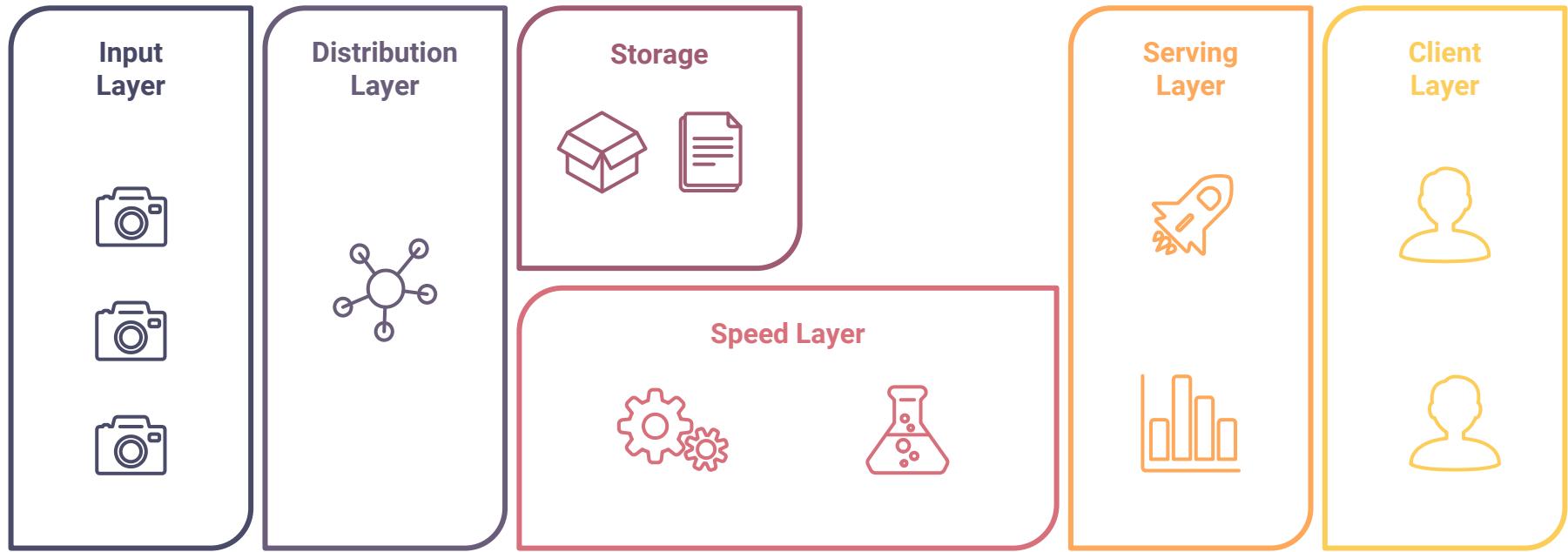
- Apache Storm

Analytics

- YOLO

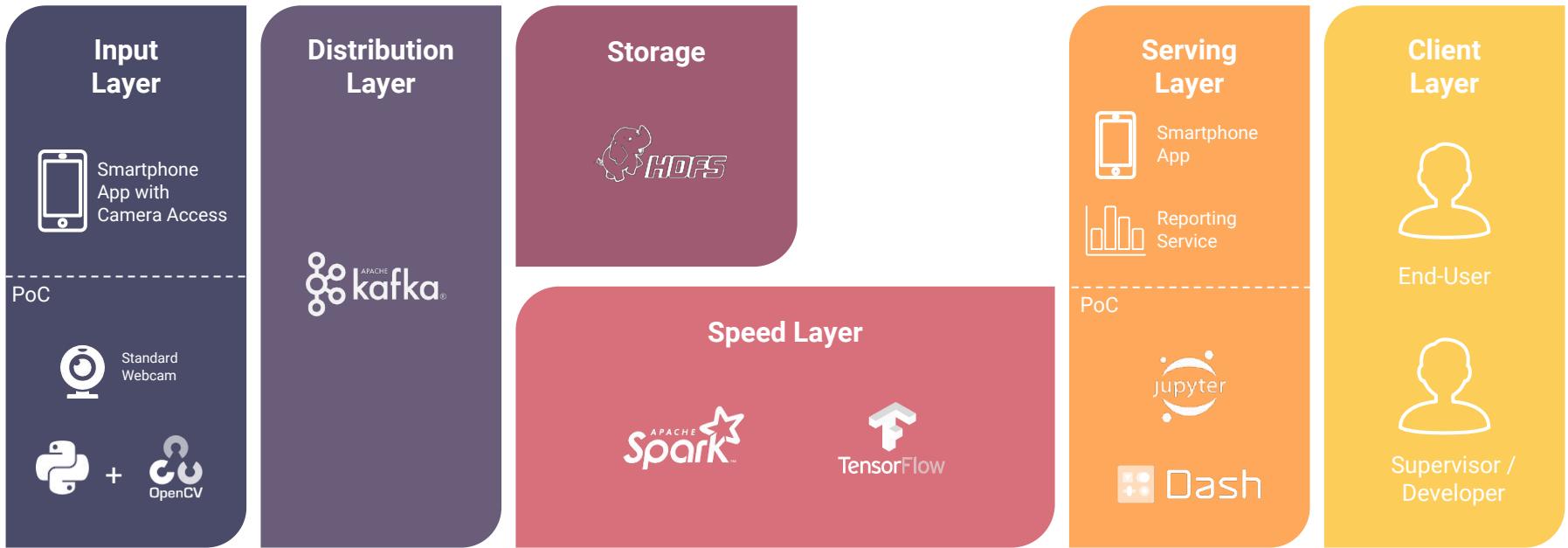
REFERENCE ARCHITECTURE

Kappa Layer Model



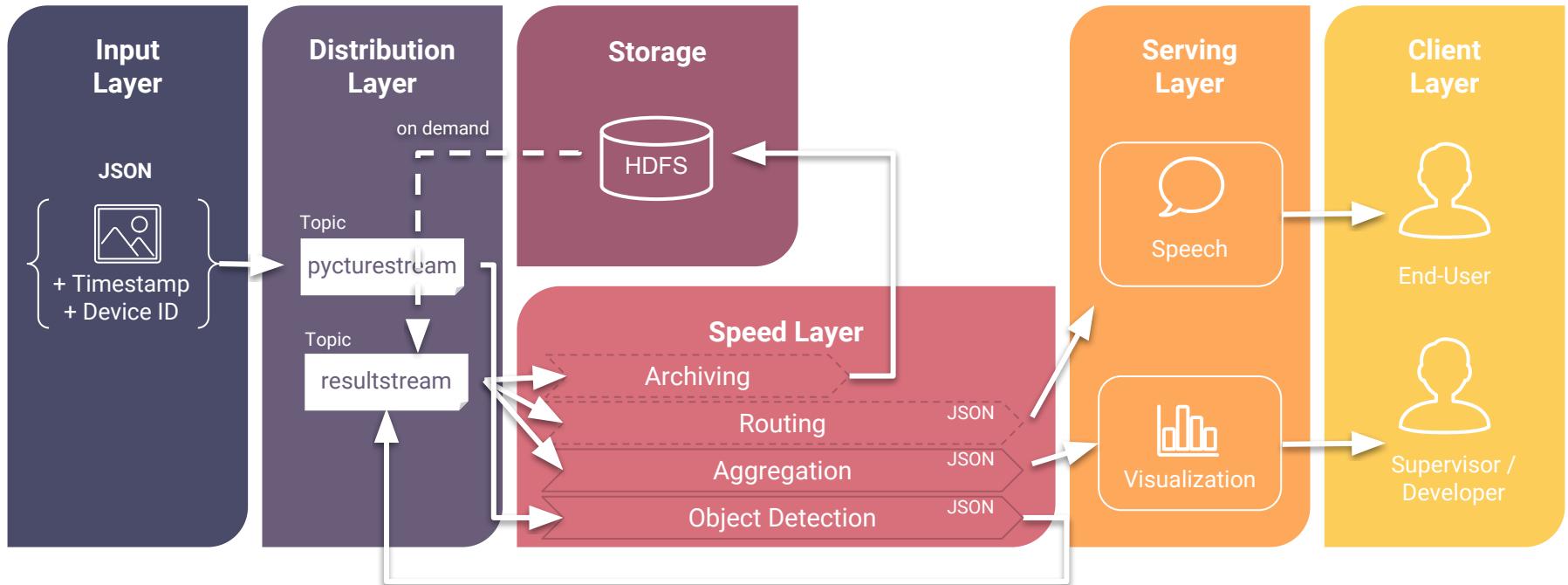
REFERENCE ARCHITECTURE

Component View



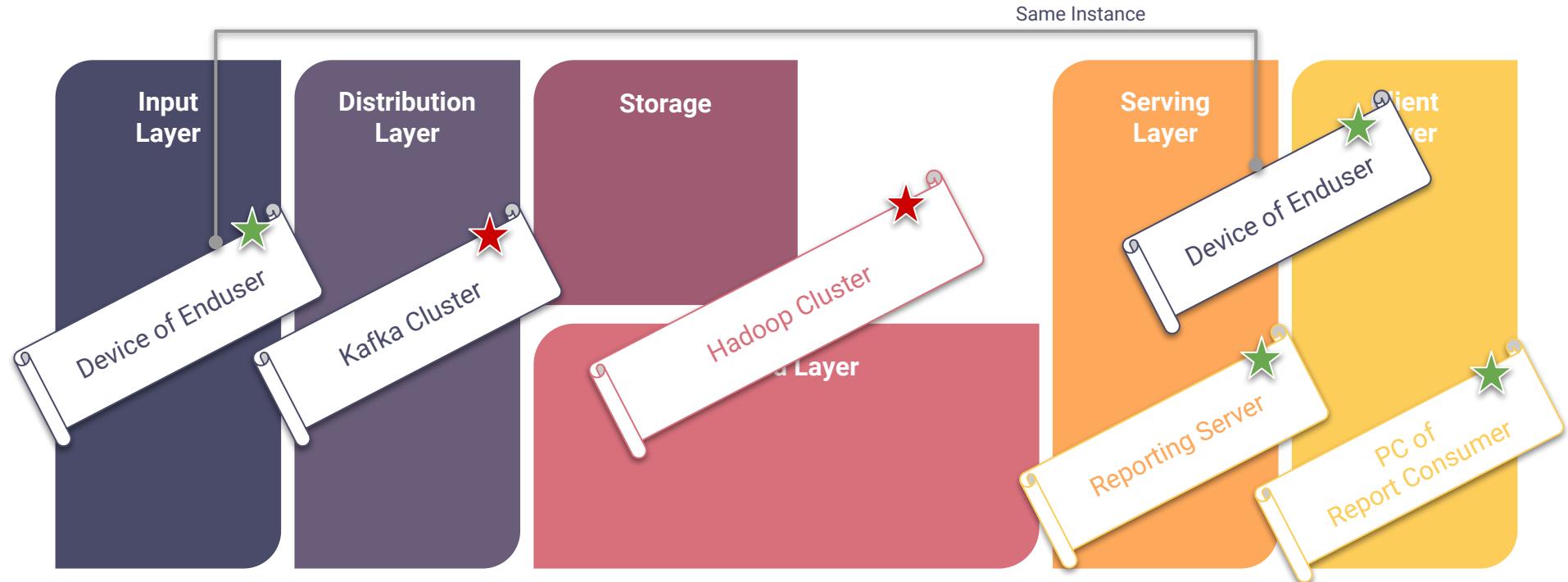
REFERENCE ARCHITECTURE

Information Flow



REFERENCE ARCHITECTURE

Deployment View



For PoC ★ Developer-Notebook

★ VM on Developer-Notebook

IMPLEMENTATION DETAILS

Client-App

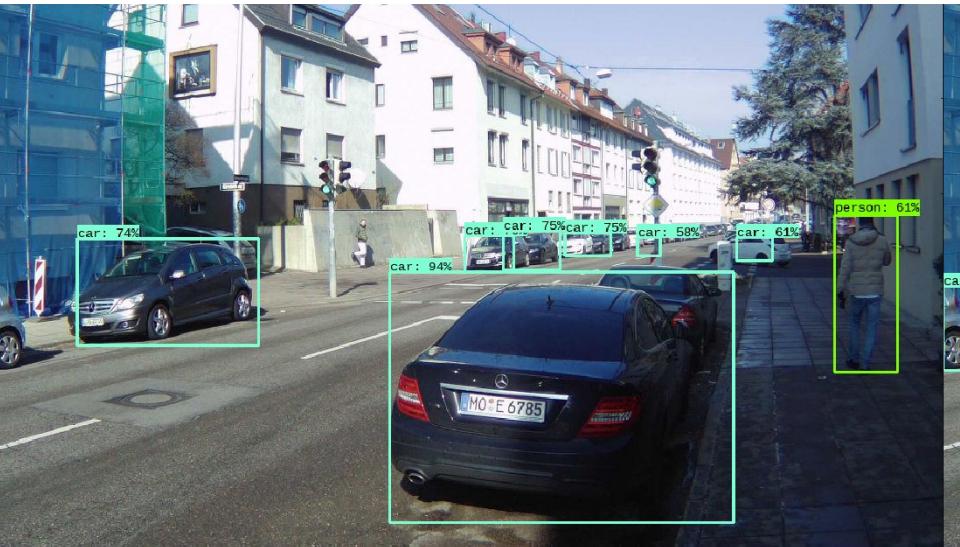
Using OpenCV to capture Webcam Image

Serialize Image and send to Kafka using kafka-python

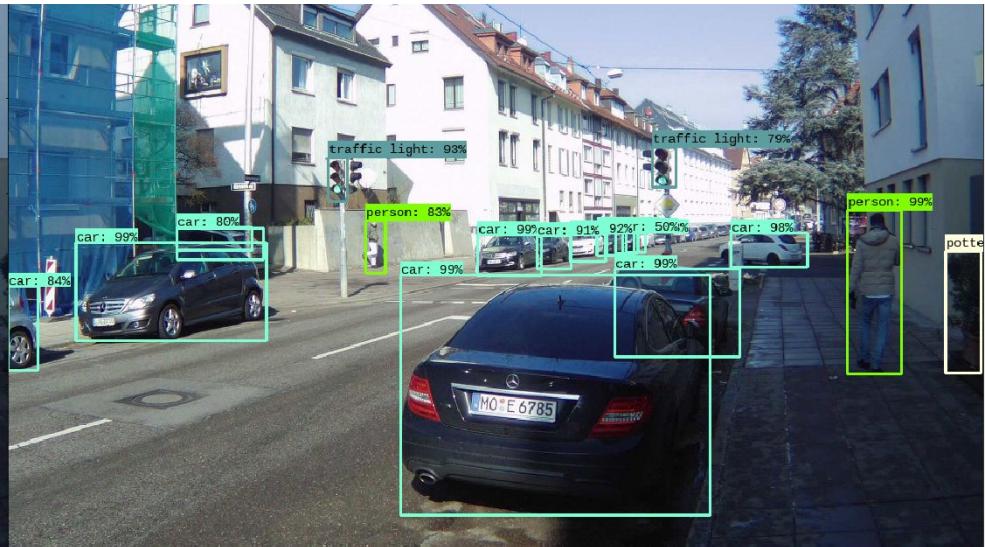
```
51      # Start Streaming...
52      logger.info('-'*50)
53      self.stream_video()
54
55  def stream_video(self):
56      """Start streaming video frames to Kafka forever."""
57      logger.info(f'Start capturing frames every {self.interval} sec.')
58      while True:
59          vidcap = cv2.VideoCapture(self.video_source)
60          vidcap.set(3,1280)
61          vidcap.set(4,720)
62          success, image = vidcap.read()
63          timestamp = dt.datetime.now().isoformat()
64          vidcap.release()
65          if success is True:
66              # Base64 encode image for transfer in json
67              jpg = cv2.imencode('.jpg', image)[1]
68              jpg_as_text = base64.b64encode(jpg).decode('utf-8')
69              # Build object and send to Kafka
70              result = {
71                  'image': jpg_as_text,
72                  'timestamp': timestamp,
73                  'camera_id': self.camera_id
74              }
75              self.send_to_kafka(result)
76              self.save_image(image)
77          else:
78              logger.error(f'Could not read image from {self.video_source}!')
79          # Sleep interval, before next capture
80          time.sleep(self.interval)
81
82  def send_to_kafka(self, data):
83      """Send JSON payload to topic in Kafka."""
84
```

IMPLEMENTATION DETAILS

Model Evaluation



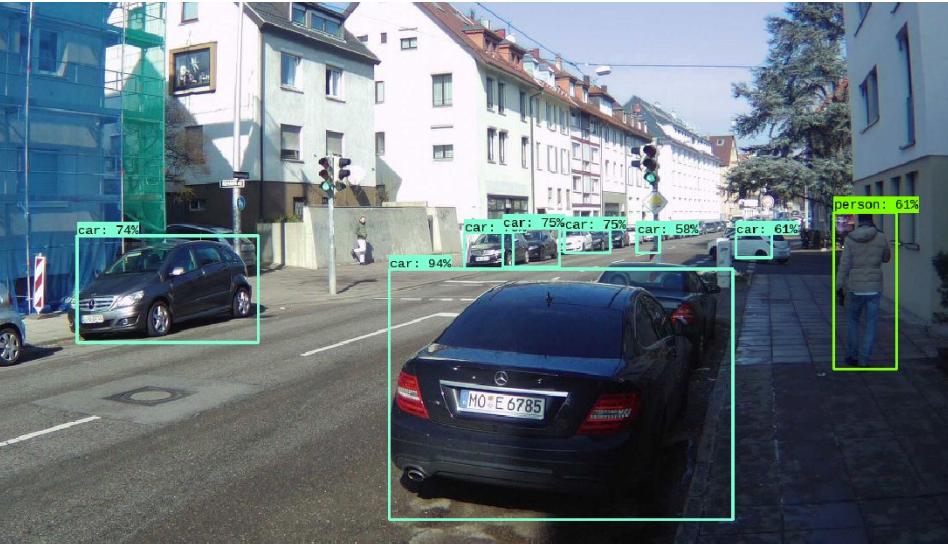
ssd_mobilenet_v1_coco (4.2 sec)



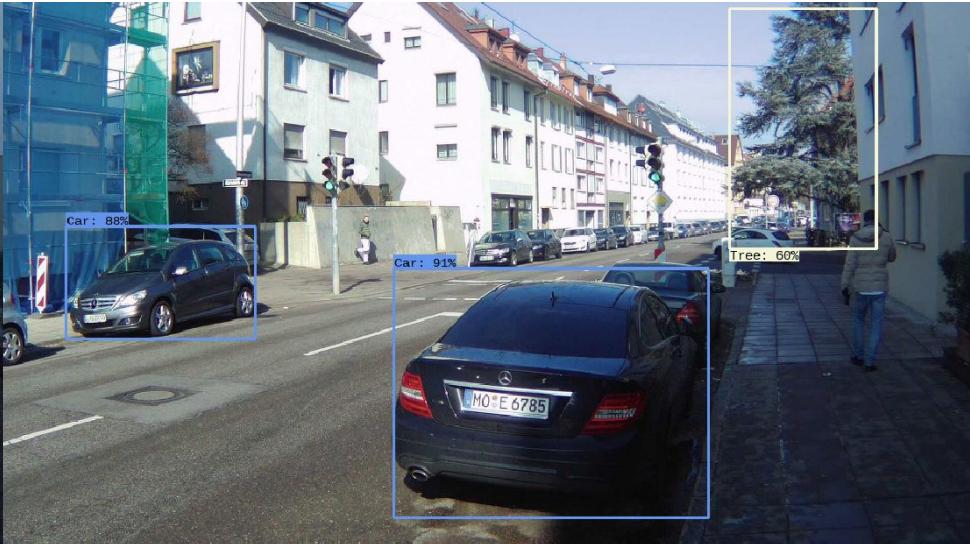
faster_rcnn_resnet101_coco (17.9 sec)

IMPLEMENTATION DETAILS

Model Evaluation



ssd_mobilenet_v1_coco (4.2 sec)



faster_rcnn_inception_resnet_v2_atrous_oid (83.7 sec)

IMPLEMENTATION DETAILS

Processing

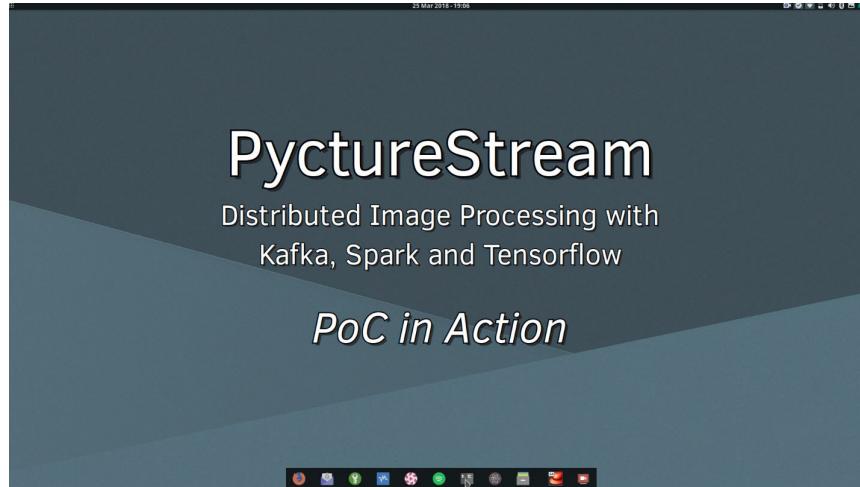
Create SparkContext and broadcast Detection Model to Worker Nodes

Listen to Kafka Topic using Spark Streaming

Run TensorFlow with Object Detection API on Images

DEMO OF THE PROTOTYPE

Demonstrating the Pipeline



Considerations



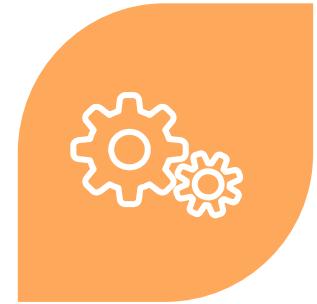
CLOUD vs. ON-PREMISE

Tough decision. Has to be decided based on individual Circumstances.



PROTOTYPE ONLY

The Demo ran on a single Machine with lots of Overhead.
Performance would scale!



STREAMS ARE DIFFERENT

Stream Processing needs a special Skill Set.
Like the one our Experts have.



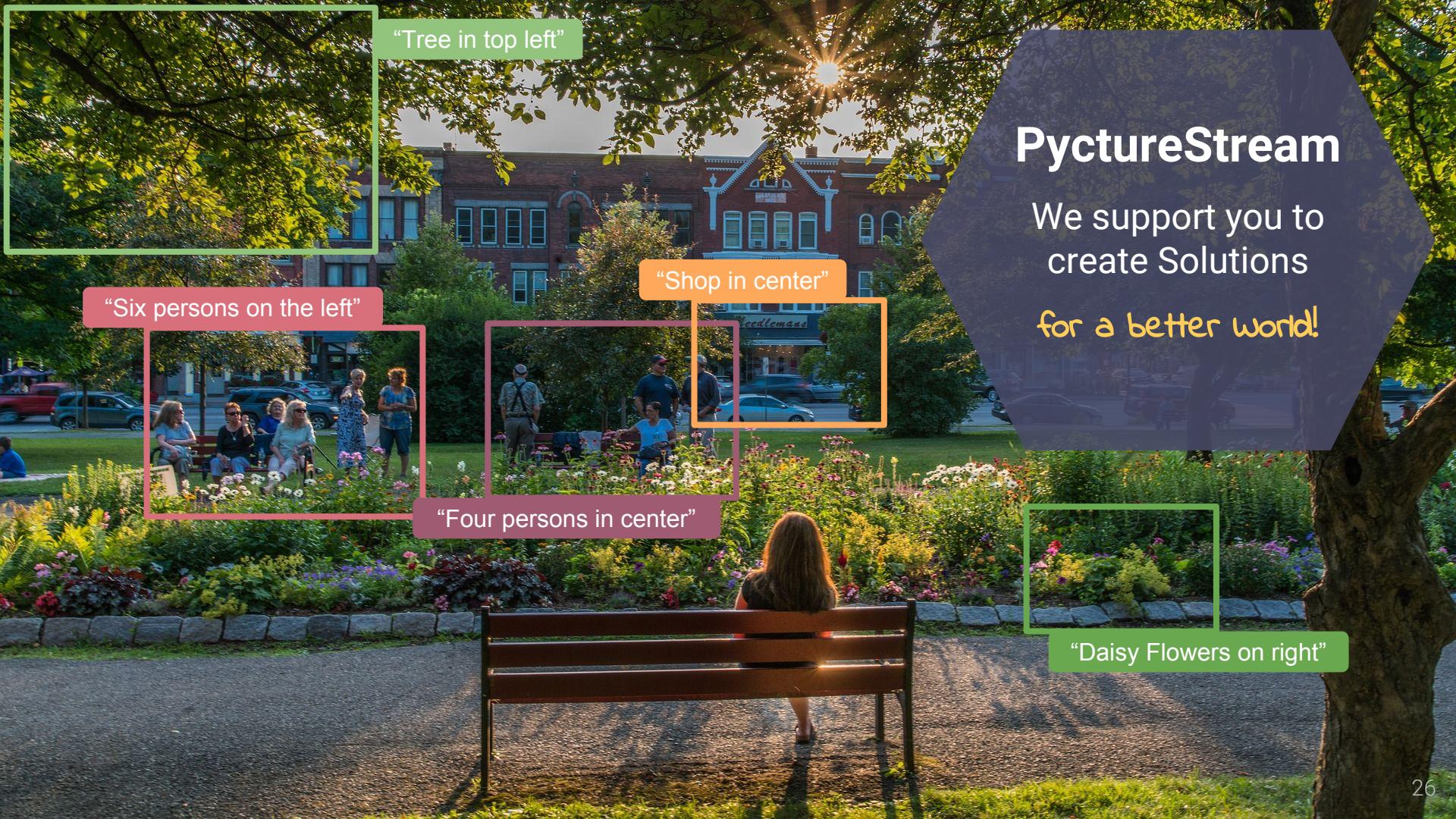
END-TO-END

We help you to implement a complete solution.
From Idea to shipped Product.



PyctureStream

We support you to
create Solutions



PyctureStream

We support you to
create Solutions

for a better world!

Q & A

Please ask your Questions now! :-)

References

Aichele C. (2006): Intelligentes Projektmanagement. Stuttgart: Kohlhammer Verlag.

Be My Eyes (2018): Be My Eyes – Hilfe für Blinde und Sehbehinderte. Android App im Google Play Store. URL: <https://play.google.com/store/apps/details?id=com.bemyeyes.bemyeyes>. Zugriff: 25.03.2018.

Büch H., alias dynobo (2018). PyctureStream – PoC for image processing using Kafka, Spark Streaming & TensorFlow. GitHub repository. URL: <https://github.com/dynobo/PyctureStream>. Zugriff: 26.02.2018.

BMI – Business Models Inc. (2018): The Business Model Canvas – Tool to help you understand a business model. URL: <https://www.businessmodelsinc.com/tools-skills/tools/business-model-canvas>. Zugriff: 25.03.2018.

Cloudera Inc. (2018). QuickStarts for CDH 5.12 – Virtualized clusters for easy installation on your desktop. URL: https://www.cloudera.com/downloads/quickstart_vms/5-12.html. Zugriff: 26.02.2018.

Dell Inc. (2017). Solution Overview – Dell EMC Ready Bundles for Hadoop. URL: <https://www.emc.com/collateral/solution-overview/ready-bundles-for-hadoop-solution-overview.pdf>. Zugriff: 13.03.2018.

Google LLC (2018): Cloud Vision API – Quotas and Limits. URL: <https://cloud.google.com/vision/quotas>. Zugriff: 22.03.2018.

Hunter T. (2016): Deep Learning with Apache Spark and TensorFlow. URL: <https://databricks.com/blog/2016/01/25/deep-learning-with-apache-spark-and-tensorflow.html>. Zugriff: 28.02.2018.

LeaseWeb (2018): Private Rack Colocation – Set your business free with powerful private racks. URL: <https://www.leaseweb.com/colocation/private-rack>. Zugriff: 20.03.2018

Pathirage M. (2017): Kappa Architecture – Where everything is a stream. URL: <http://www.kappa-architecture.com>. Zugriff: 23.03.2018.

References

- Mallot** H. A. (2006): Visuelle Wahrnehmung. In J. Funke & P. A. Frensch (Hrsg.): Handbuch der Allgemeinen Psychologie – Kognition. Göttingen: Hogrefe Verlag.
- Moffat** R. (2017): Getting Started with Spark Streaming, Python and Kafka. URL:
<https://www.rittmanmead.com/blog/2017/01/getting-started-with-spark-streaming-with-python-and-kafka/>. Zugriff: 05.03.2018.
- Plotly** Inc. (2018): Dash by Plotly – Build beautiful web-based interfaces in Python. URL: <https://plot.ly/products/dash/>. Zugriff: 20.03.2018.
- Project Jupyter** (2018, 20. Februar): JupyterLab is Ready for Users. *Jupyter Blog – All the latest about Project Jupyter*. URL:
<https://blog.jupyter.org/jupyterlab-is-ready-for-users-5a6f039b8906>. Zugriff: 25.02.2018.
- Oracle** Corp. (2018). Oracle VM VirtualBox. URL: <https://www.virtualbox.org>. Zugriff: 26.02.2018.
- TensorFlow** (2018a): Tensorflow Object Detection API. GitHub repository. URL: https://github.com/tensorflow/models/tree/master/research/object_detection. Zugriff: 26.02.2018.
- TensorFlow** (2018b): Tensorflow detection model zoo. GitHub repository. URL:
https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md. Zugriff: 26.02.2018.

Image Sources

Photographs are from pixabay.com and licensed CC0
Icons are from SlidesCarnival.com and licensed CC BY 4.0
Charts & Screenshots are self-made

BACKUP SLIDES

If enough is not enough...

Cost Calculation Cloud

Compute Engine
1 x
730 total hours per month
VM class: regular
Instance type: n1-standard-4
Region: Frankfurt
Sustained Use Discount: 30% ?
Effective Hourly Rate: \$0.171
Estimated Component Cost: \$125.09 per 1 month
Persistent Disk
Frankfurt
Storage: 200 GB
\$9.60
Cloud Dataflow
3 x n1-standard-4 workers in Streaming Mode
Region: Europe
Total vCPU Hours: 8,640
Total Memory Hours: 32,400 GB/Hours
PD Local Storage: 64,800 GB/hours
\$760.75



Cloud Machine Learning
Region: United States
ML Training Units: 2.473
Job run time: 0 minutes
Prediction Mode: online
of Predictions: 2,591,000
Total Inference time: 5,182,000 seconds
\$431.83
Total Estimated Cost: \$1,327.28 per 1 month
Cloud Dataproc
Cluster size: 176 vCPUs
Time used: 720 hours
\$1,267.20
Total Estimated Cost: \$2,162.64 per 1 month

Cloud Dataproc
Cluster size: 176 vCPUs
Time used: 720 hours
\$1,267.20
Total Estimated Cost: \$2,162.64 per 1 month

OR

Vision API
Label Detection: 864,000
OCR: 0
Explicit Content Detection: 0
Facial Detection: 0
Landmark Detection: 0
Logo Detection: 0
Image Properties: 0
Web Detection: 0
Document Text Detection: 0
\$1,294.50
Total Estimated Cost: \$2,189.94 per 1 month

Cost Calculation On Premise

The screenshot shows the Dell website's product page for the PowerEdge R730xd Rack Server. At the top, there's a navigation bar with links for Products, Solutions & Services, Support, and Deals. A search bar and a sign-in link are also present. Below the navigation, a promotional banner offers a free dock with select Latitude laptops or Precision workstations, or a 10% discount on Inspiron and XPS PCs with coupon SAVE10. A "Shop Now" button and a "Click to Chat" link are included. The main content area is titled "PowerEdge R730xd Rack Server Summary". It displays the Dell Price as \$4,187.45. To the right of the price is an "Add to Cart" button. Below the price, it shows starting at \$6,690.00, total savings of \$2,502.55, and standard delivery is free. There's also a section for Dell Business Credit with a note about low rates and an "Apply" button, along with a "View Delivery Dates" link.

- 8 Server for Hadoop-Cluster
(3 Name Nodes, 1 Edge Node, 4 Worker Nodes)
- 3 Server for Kafka-Cluster
(1 Zookeeper Node, 2 Broker Nodes)
- Reporting will run on Edge Node
- 700 \$ Rent per month for Data Center
- No Licence Cost (Apache Distro)