

# Directions in Machine Learning

Michele Covell  
Google Research

with many thanks to Rahul Sunkthakar and  
many Google Machine Perception Teams  
for the slide materials



# Directions in ML

Need to understand **intent** and **importance** in audio and video

- Major progress on: detection, categorization, embeddings, 3D modeling, people-centric information, action/interaction recognition
- missing piece: need to determine **what to leave out**
  - want a **synopsis** for **authored** (and situated) media
  - want less constrained interactions for **live** situations

Generating new **creative** content can help highlight shortcomings  
(as well as providing useful content)

# Understanding: Starting from...

## Classification

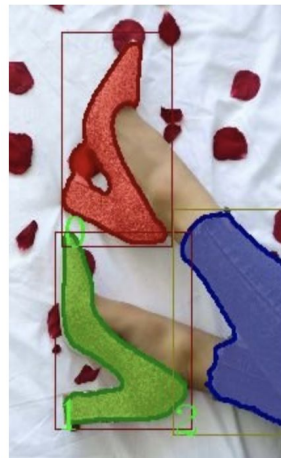
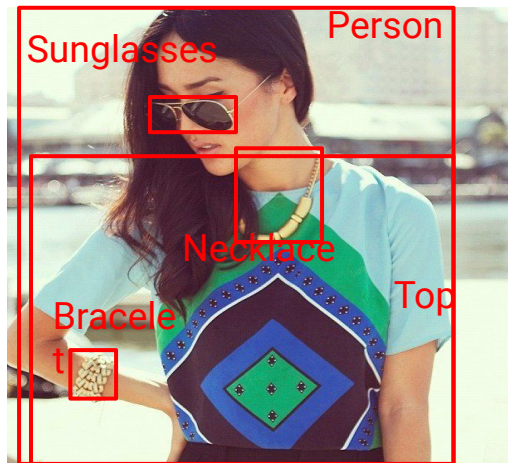
Images  $\Rightarrow$  Labels



Cake (0.93669)  
Dessert (0.91911)  
Birthday (0.89697)  
Child (0.89183)  
Fondant (0.88340)  
Birthday cake (0.87525)

## Detection

Images  $\Rightarrow$  Labelled Boxes or Regions



## Embedding

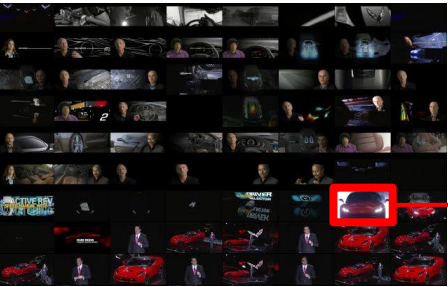
Images  $\Rightarrow$  Features



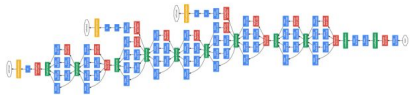
*Similar Images and Labels*



# Deep-learned visual features



Videos with **noisy** labels

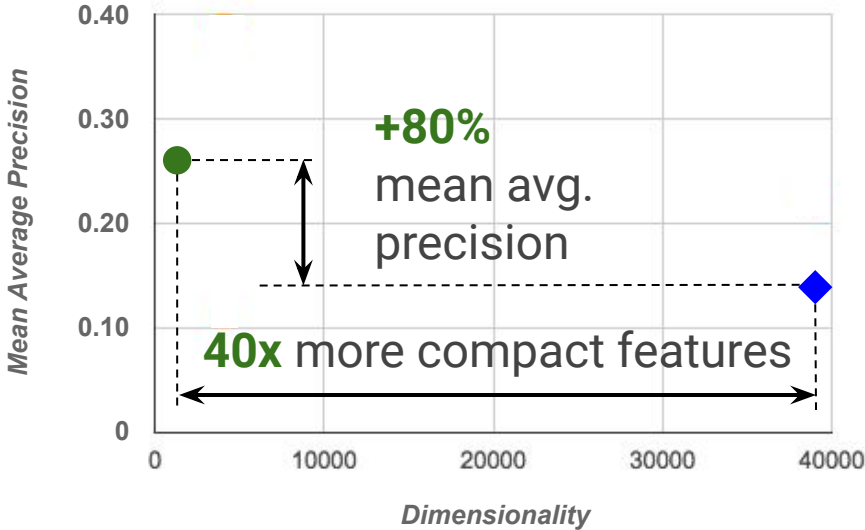


Inception model  
trained on noisy  
data (images)

Frame-level  
Bottleneck  
embedding  
layer (1000-d)

Video-level  
- Max pooling  
- Avg pooling  
- VLAD pooling

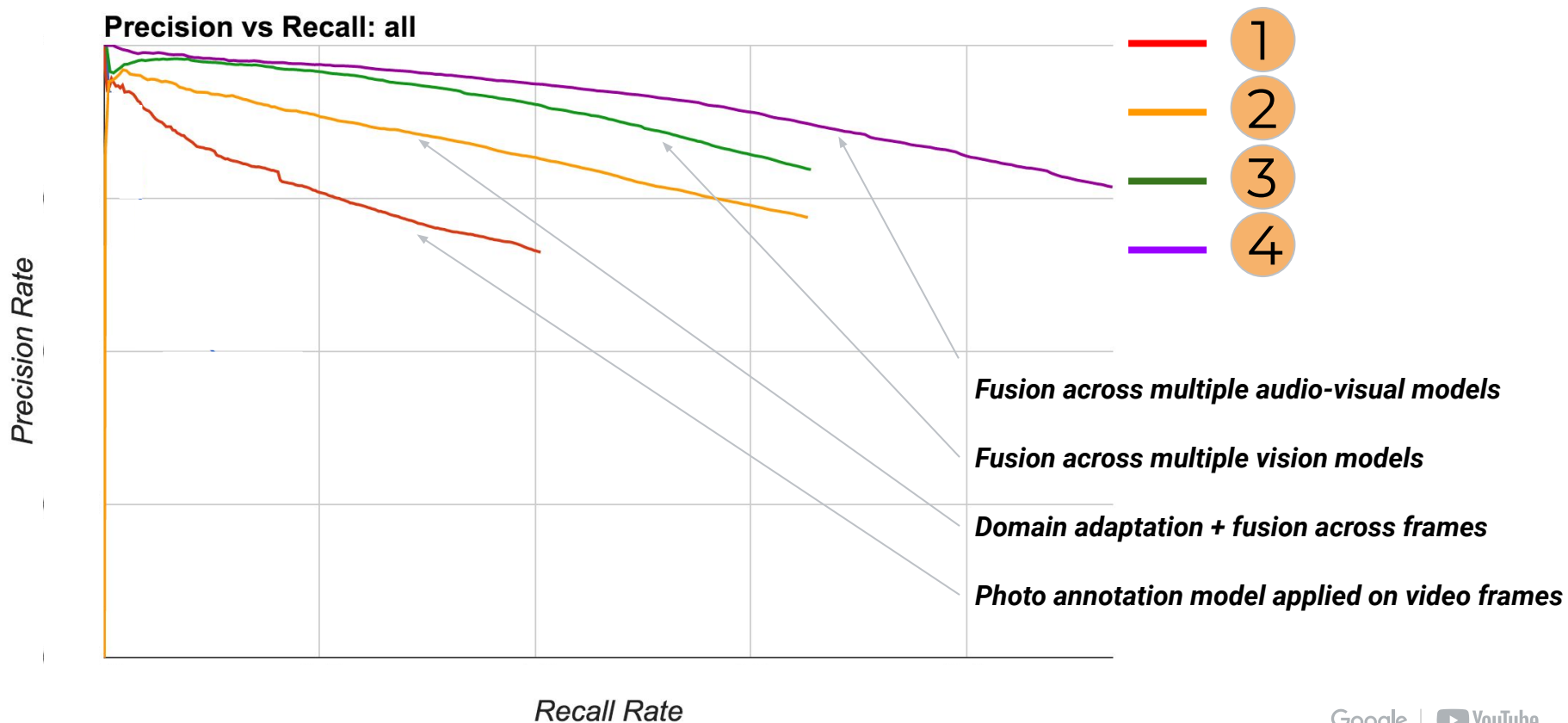
### Comparison of Features on Video Annotation



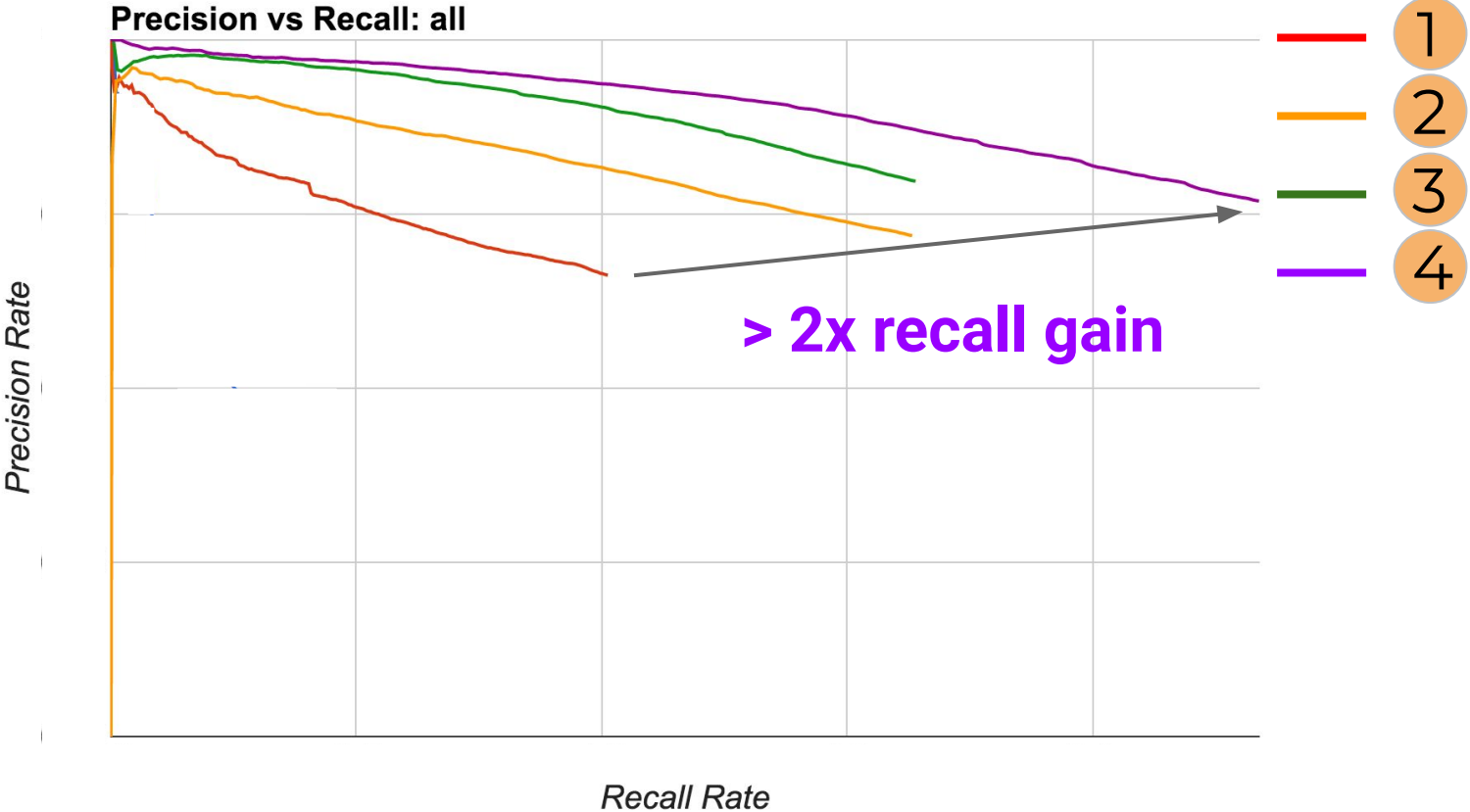
- Deep learned visual features, VLAD coding: 1024-d, 0.272 MAP
- ◆ Handcrafted audio-visual features: ~40K-d, 0.153 MAP



# Improving detection in video (starting from images)

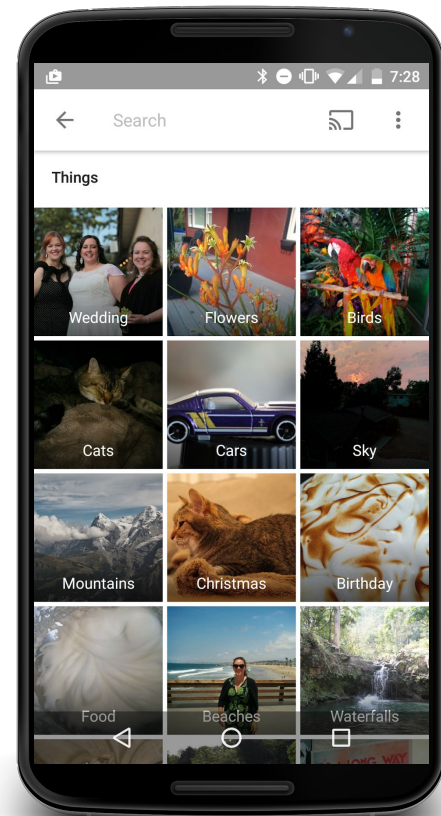
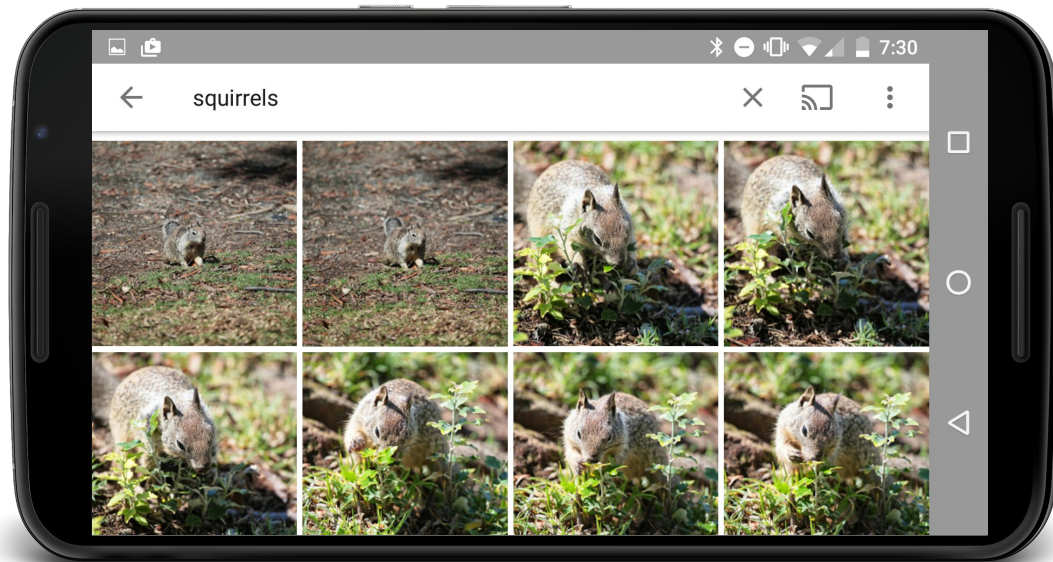


# Improving detection in video (starting from images)



# Images/Videos $\Rightarrow$ Labels

[Veit, Alldrin, Chechik, et al.]



**Automatic tagging and search!**

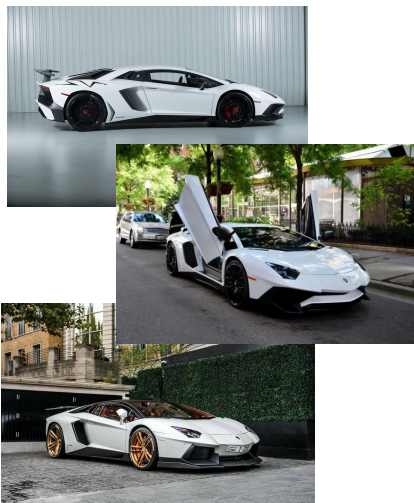
# Spectrum of semantic similarity

[Juan, Lu, Li, et al.]



Instance -> fine-grained -> coarse-grained

2016 lamborghini aventador white



lamborghini aventador red



lamborghini aventador camouflage



lamborghini huracan



2010 lamborghini Gallardo



super car

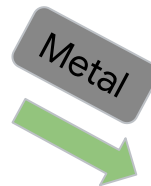
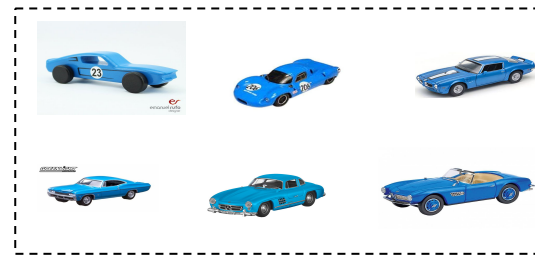
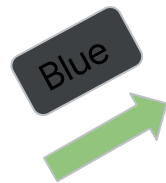
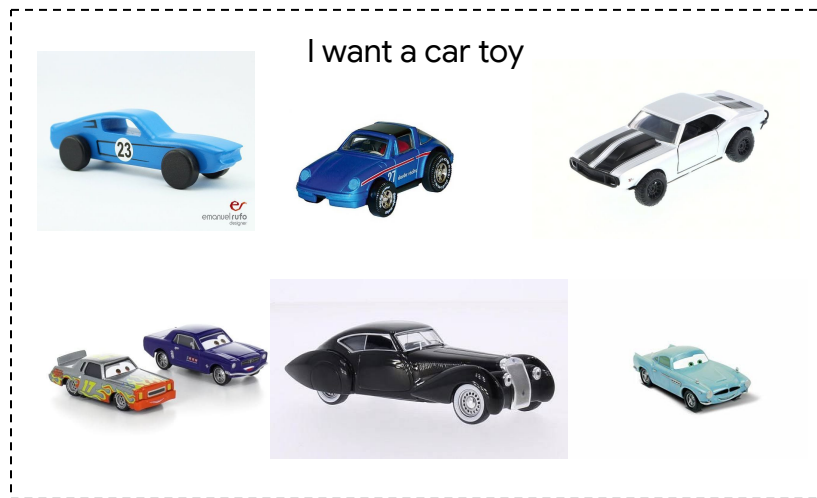


# Disentangled visual similarity

## Disentangled visual attributes.

The model is able to rank/filter images by different attribute dimensions.

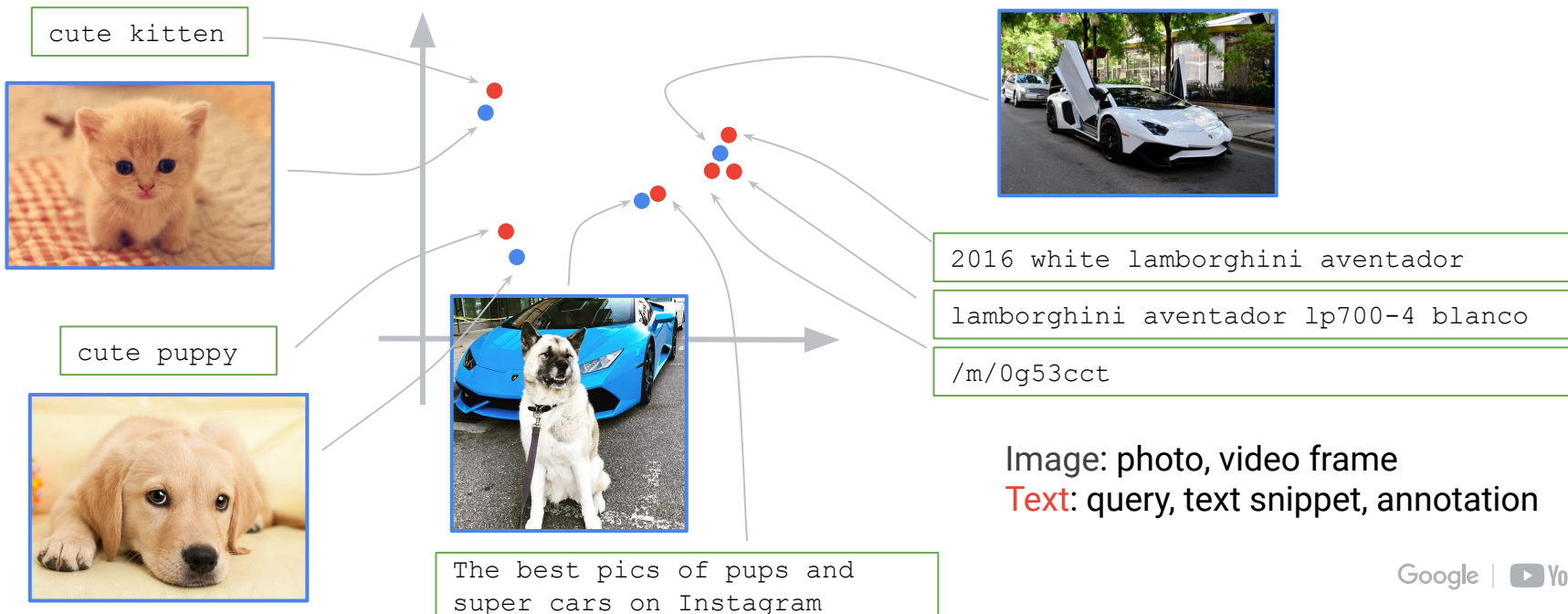
**Entangled**  $\xrightarrow{\text{disentangled attribute dimensions}}$  **Visually Fine-grained**



# Image-text co-embedding

Bridge the visual and text domain with multitask, multimodal learning

- **Multimodal:** Knowledge in one domain helps learning in the other
- **Multitask:** vast amount of visual-text and NLP data at Google
- Generic models to be used / fine-tuned for cross-domain inference / learning cases





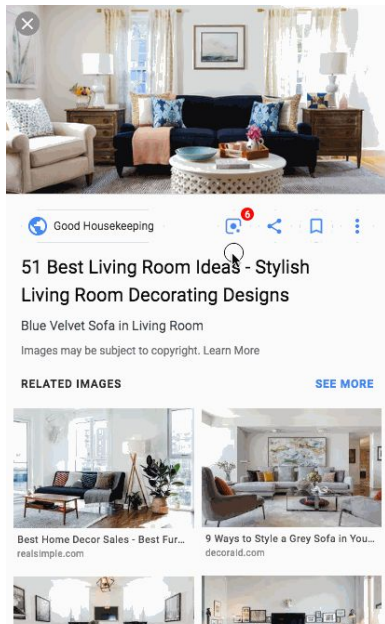
# Specialized embeddings

(e.g., products)

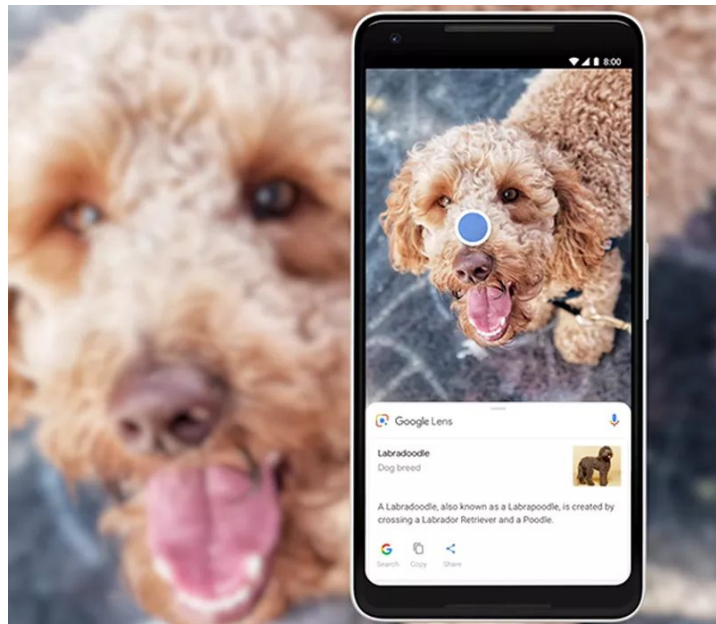


<https://cloud.google.com/vision/product-search/docs/>

# Use cases in Search and Lens



Lens in Google Images



Lens natural world, similar products, gleaming, similar images, among other improvements





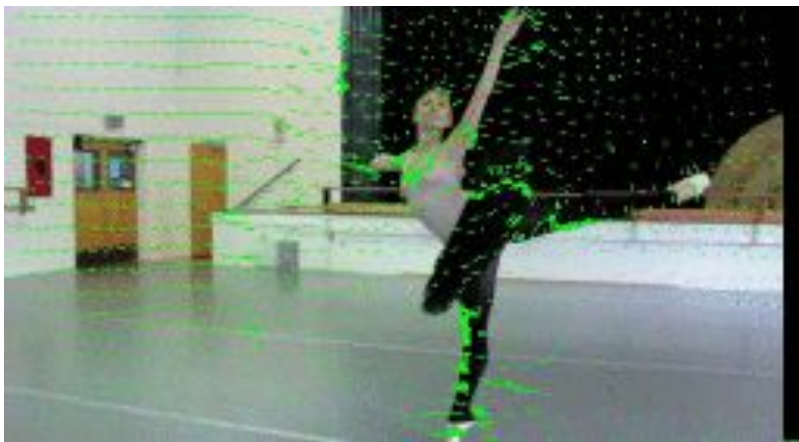
# Using Video to help with Object Recognition in Images

Weak, noisy Labels from:

- meta-data, comments
- labeling using image-trained networks



# Weakly & Self-Supervised Learning from Video



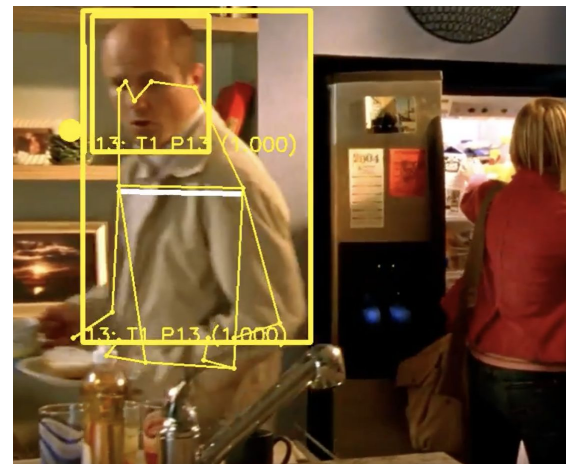
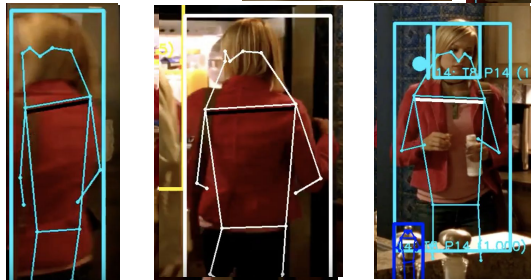
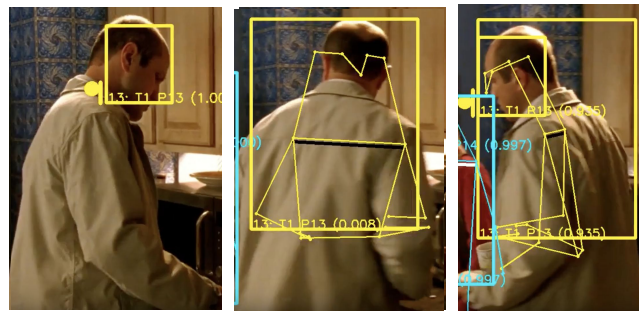
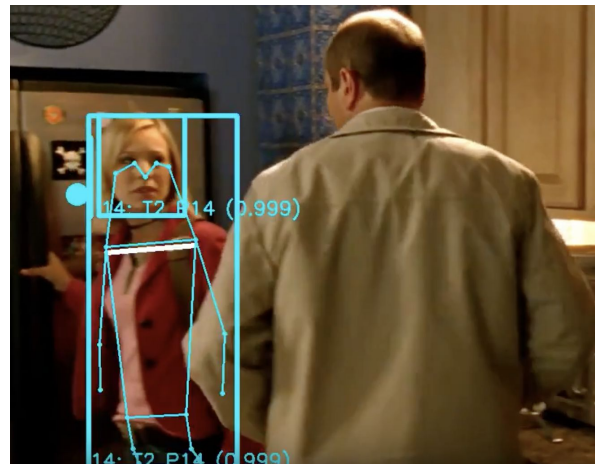
Candidate Set of Hard Positives:  
Semantically consistent frames  
from dense trajectory tracking



# Using consistency signals for supervision

**Goal:** Pose-insensitive person embedding (i.e., PersonNet)

**Solution:** 360 degree pose samples from large image / video corpus + tracking + clustering + user feedback signals



video + tracking → 360° pose training samples



# Understanding: Starting from...

## Classification

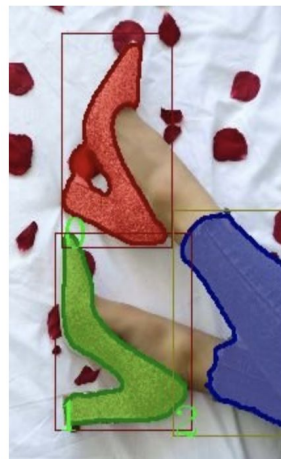
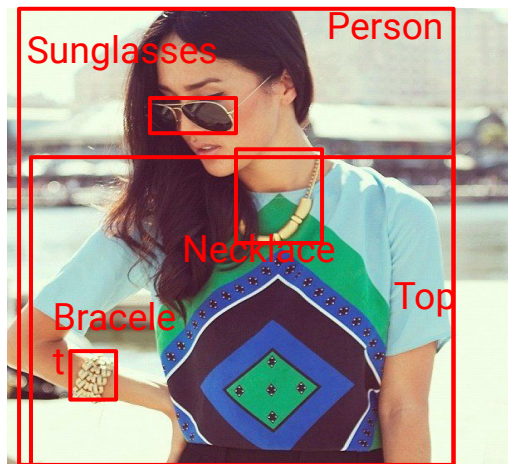
Images  $\Rightarrow$  Labels



Cake (0.93669)  
Dessert (0.91911)  
Birthday (0.89697)  
Child (0.89183)  
Fondant (0.88340)  
Birthday cake (0.87525)

## Detection

Images  $\Rightarrow$  Labelled Boxes or Regions



## Embedding

Images  $\Rightarrow$  Features



*Similar Images and Labels*

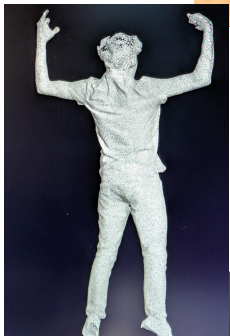


# Understanding: ... adding in...

## 3D perception

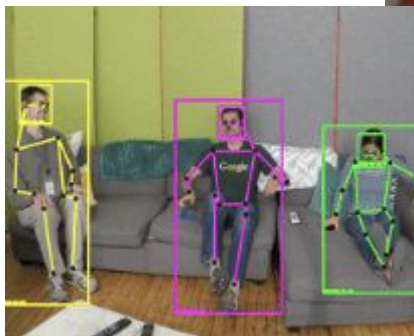
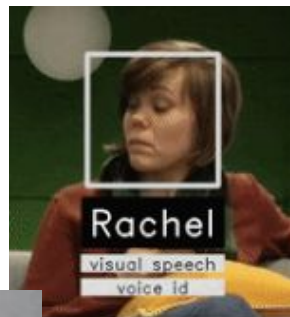
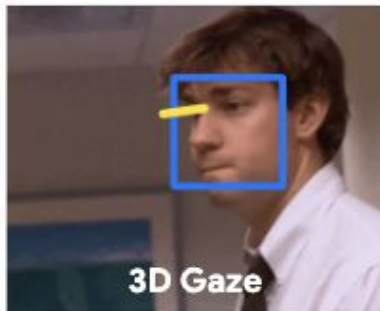
Images  $\Rightarrow$  3D relationships

single image



## Person-centric Models

Focus, Pose, Speaker Models



## Action/Interaction Recognition



YouTube  $\Rightarrow$  3D models of scenes & people

What 3D structure can we learn from watching internet video?

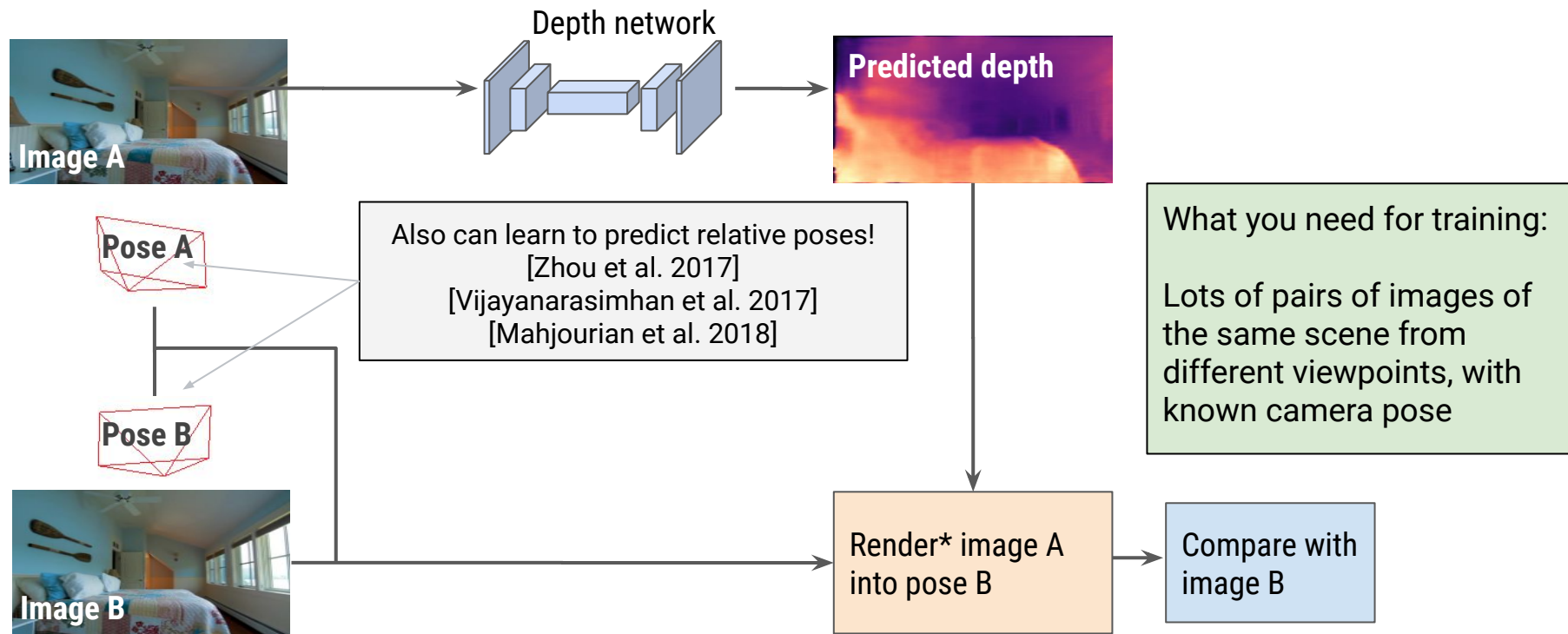
Training: Multiple views



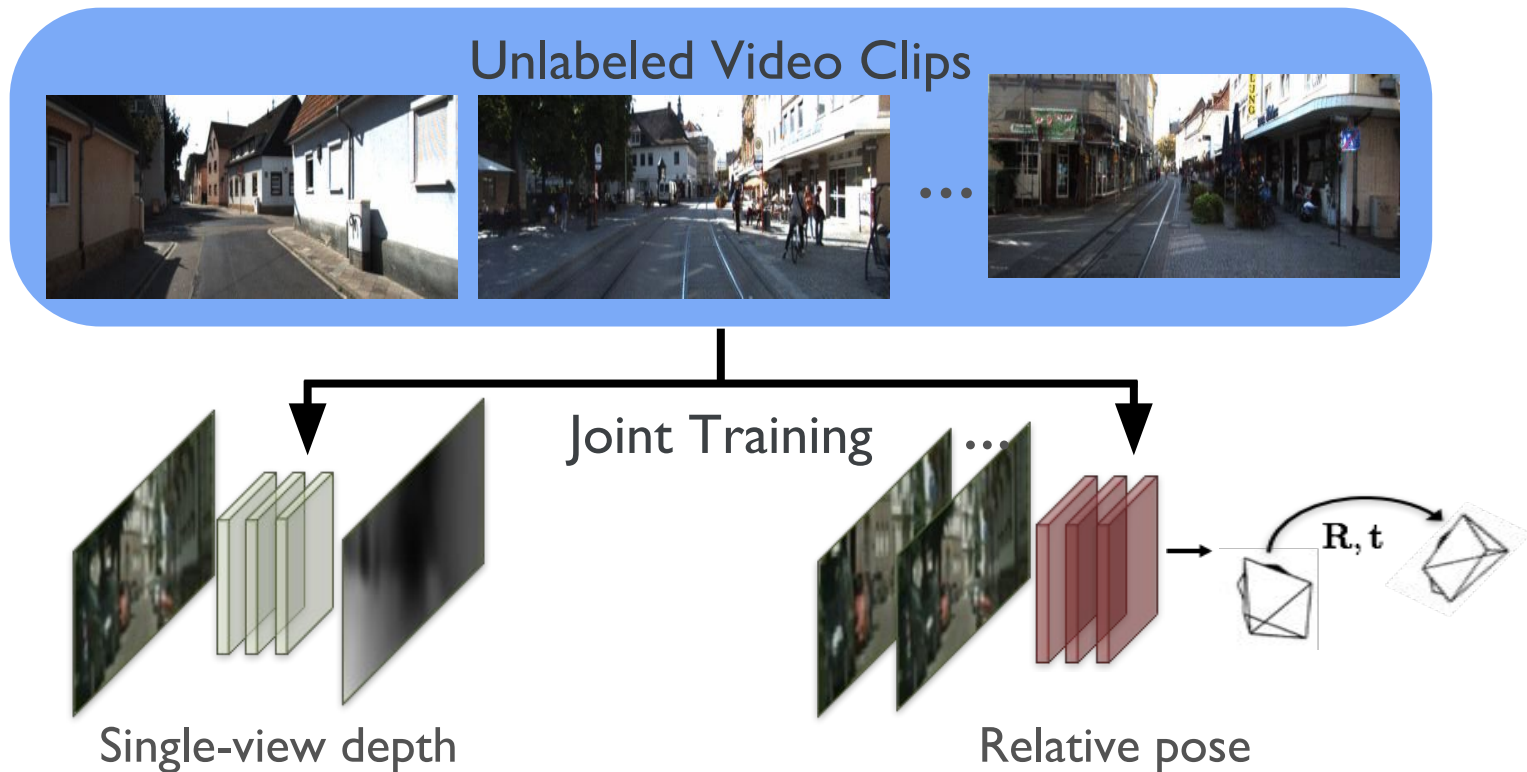
Testing: Single Image



# Beyond direct supervision

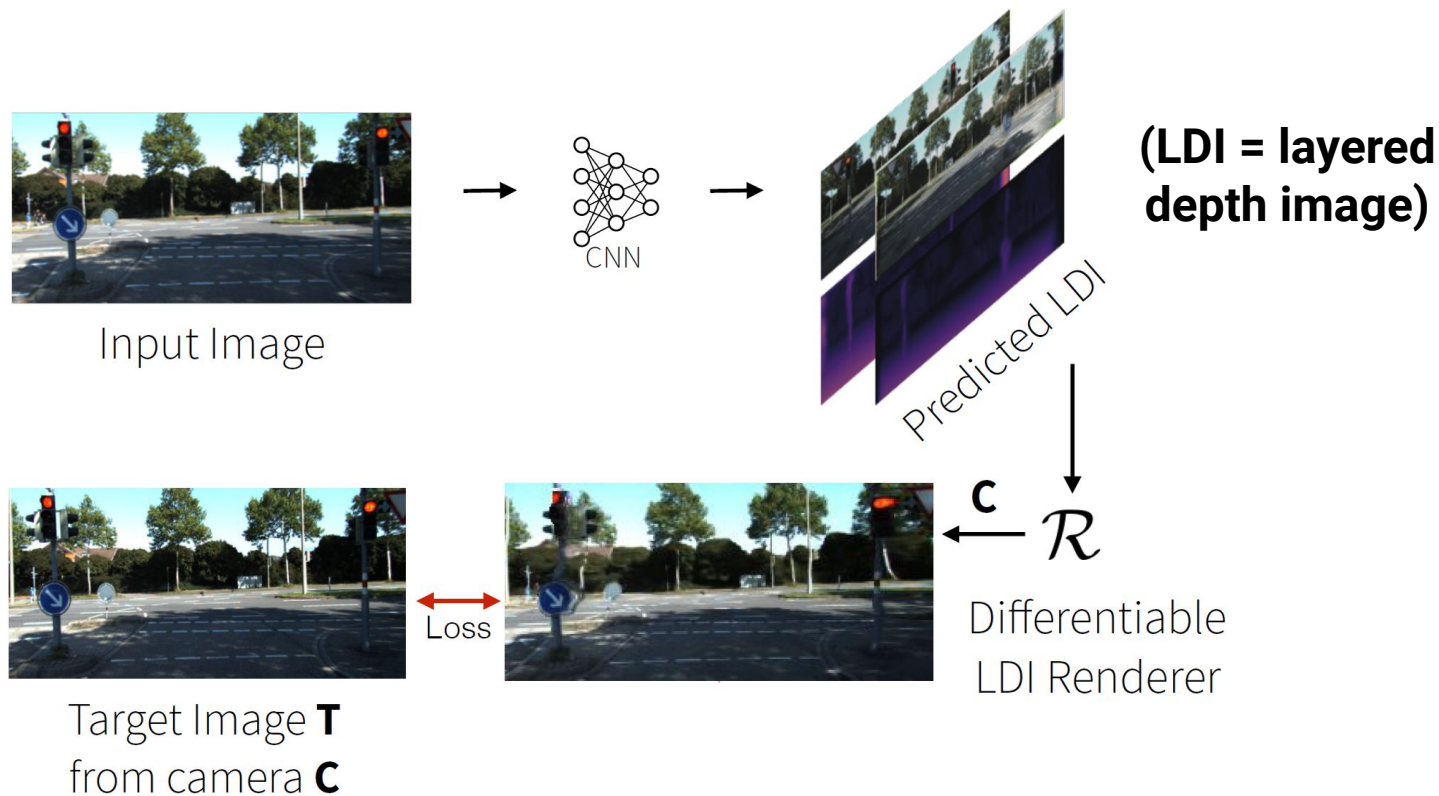


# Learning depth and camera pose via view synthesis





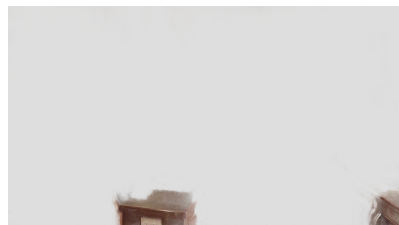
# Beyond depth maps: Learning layered models



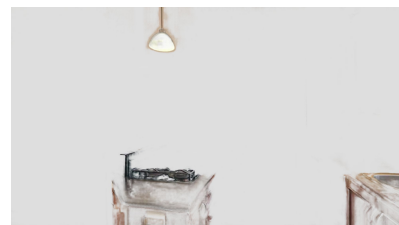
## Reference input view



Plane 0



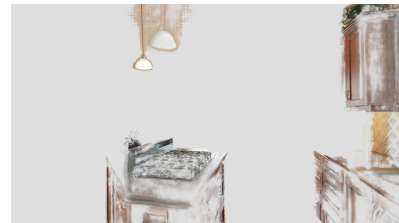
Plane 9



Plane 13



Plane 16



Plane 24



Plane 26



Output



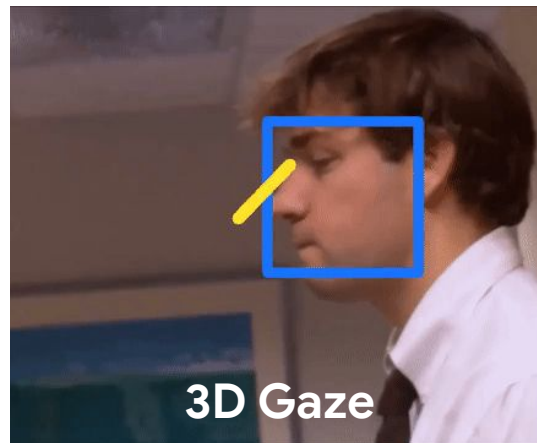




# Output



# Focus on people



- Better binary and 3D gaze models

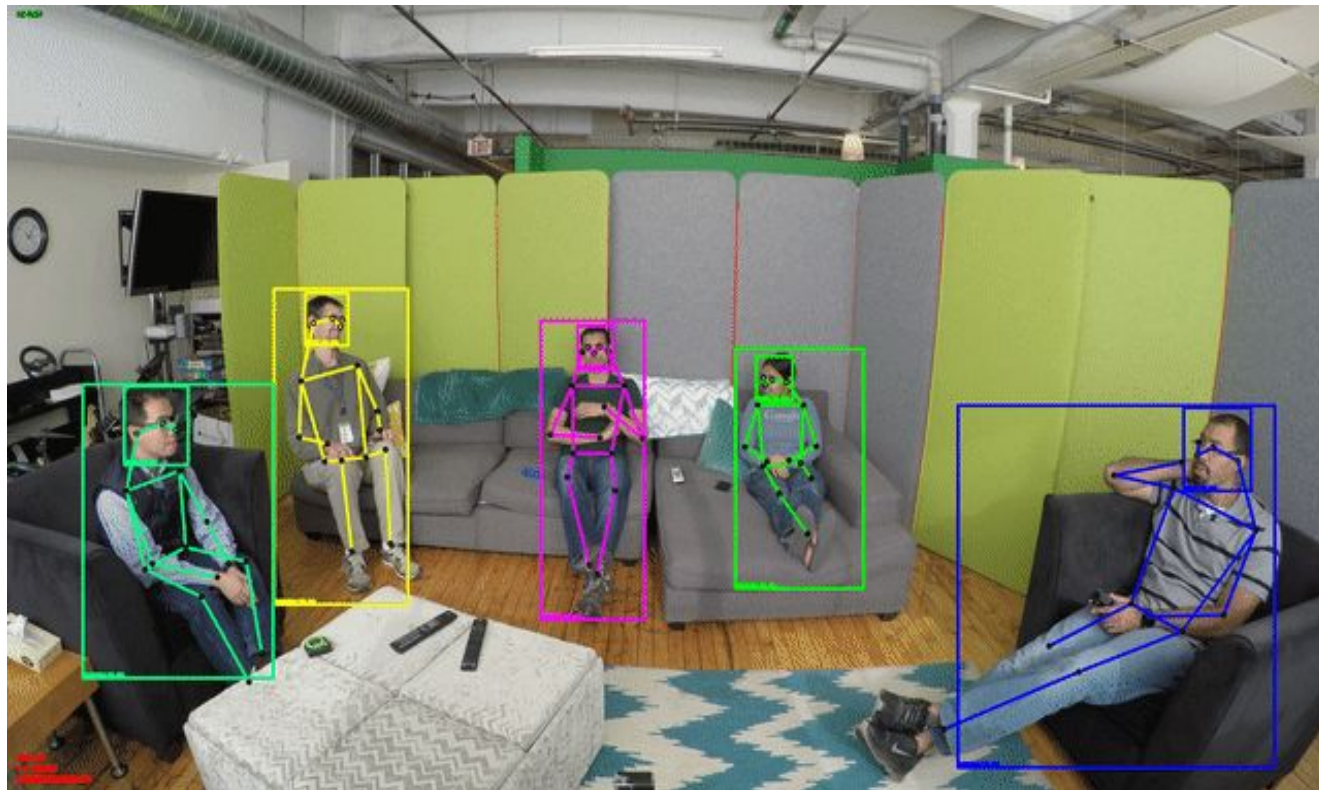
- Dynamic gestures

- Active speaker detection

- Speech Detection and Diarization



# Focus on people

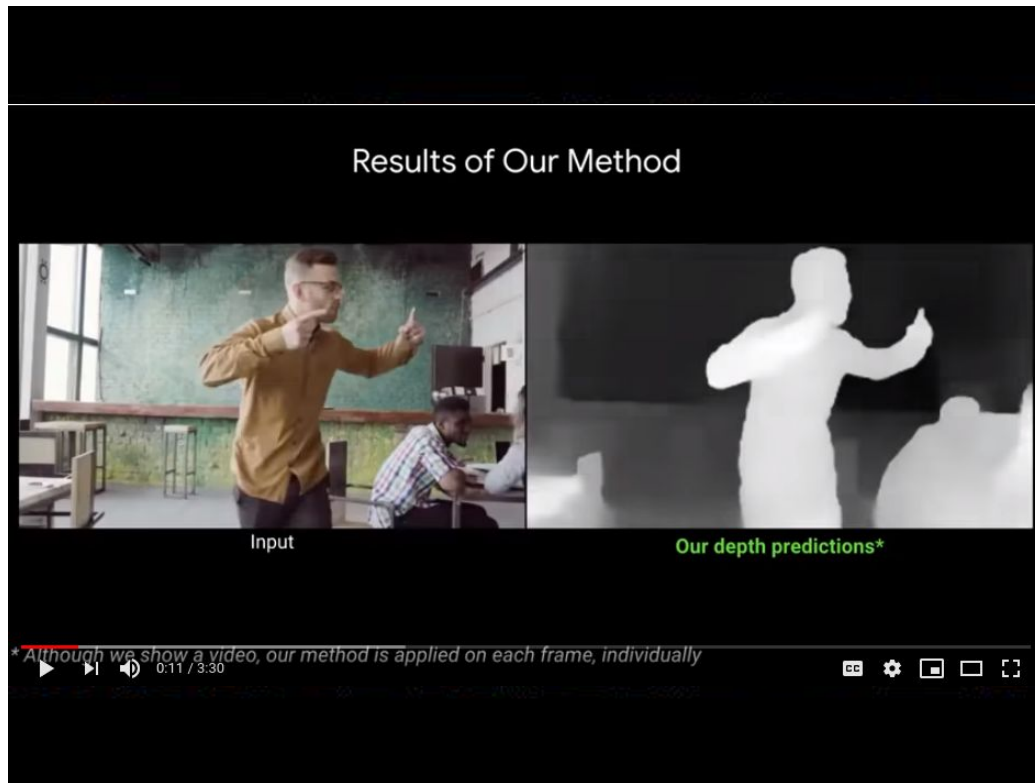


- Combined person & face SSD detection models
- Recurrent models for detection
- Probabilistic tracking
- Rotation invariance

# Depth Estimation for Moving People

Z. Li, T. Dekel, F. Cole, R. Tucker,  
N. Snavely, C. Liu, W.T. Freeman,  
CVPR 2019 Honorable Mention

Goal: depth estimation with moving camera and moving people





# Depth Estimation for Moving People


Z. Li, T. Dekel, F. Cole, R. Tucker,  
N. Snavely, C. Liu, W.T. Freeman,  
CVPR 2019 Honorable Mention

Idea: leverage MannequinChallenge dataset of **frozen** people!

Approach:  
Learn the depths of **moving people** by watching **frozen people**

**MannequinChallenge Dataset:**

- 2000 YouTube Videos
- People frozen while camera is moving
- Diverse scenes, natural human poses

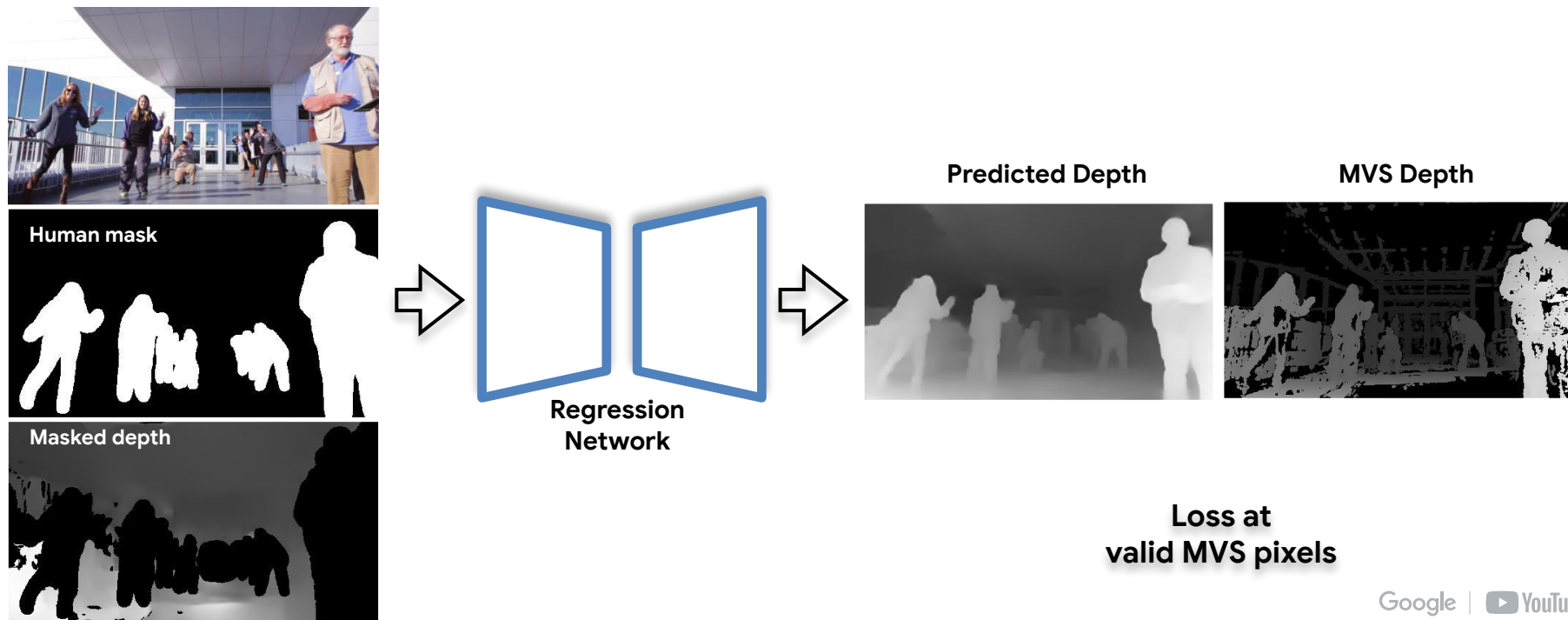


0:30 / 3:30

Google | YouTube

# Depth Estimation for Moving People

Approach: compute depth for static scene with multiview stereo, and predict depth for moving people with a regression network

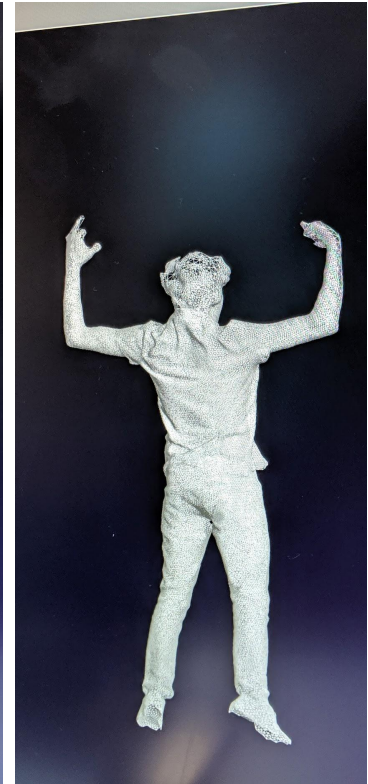
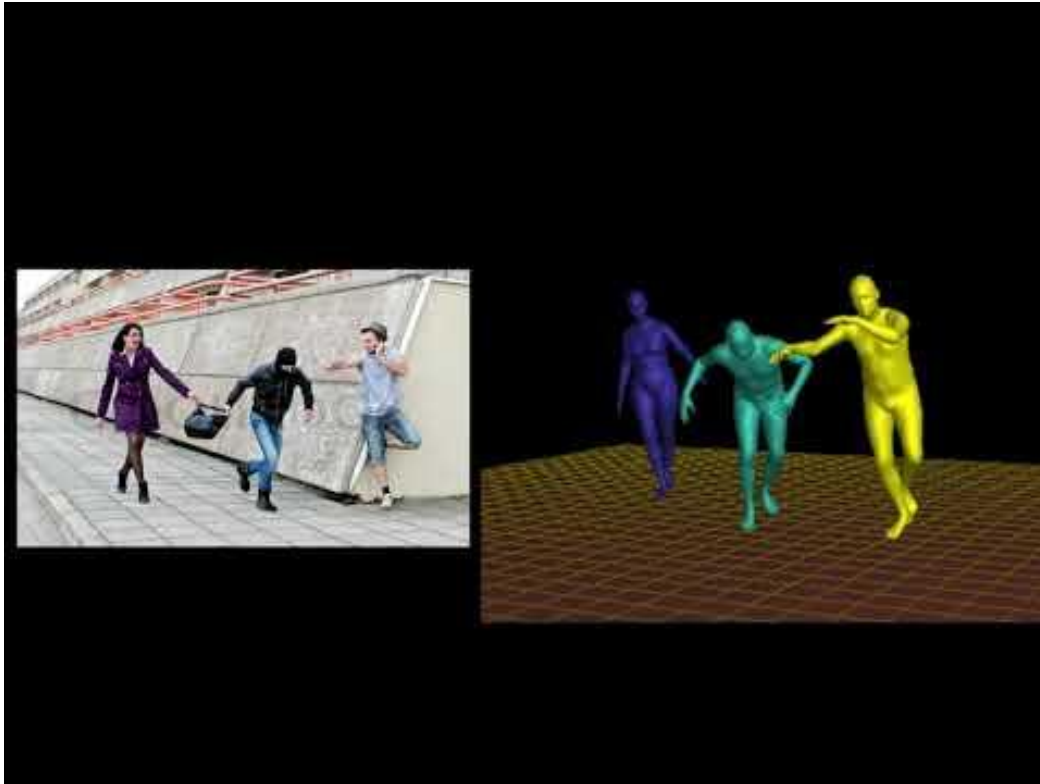


# Depth Estimation for Moving People

Result: depth estimation with moving camera and **moving** people



# Focus on People: 3D Shape Models



Physical relationships between people in 3D

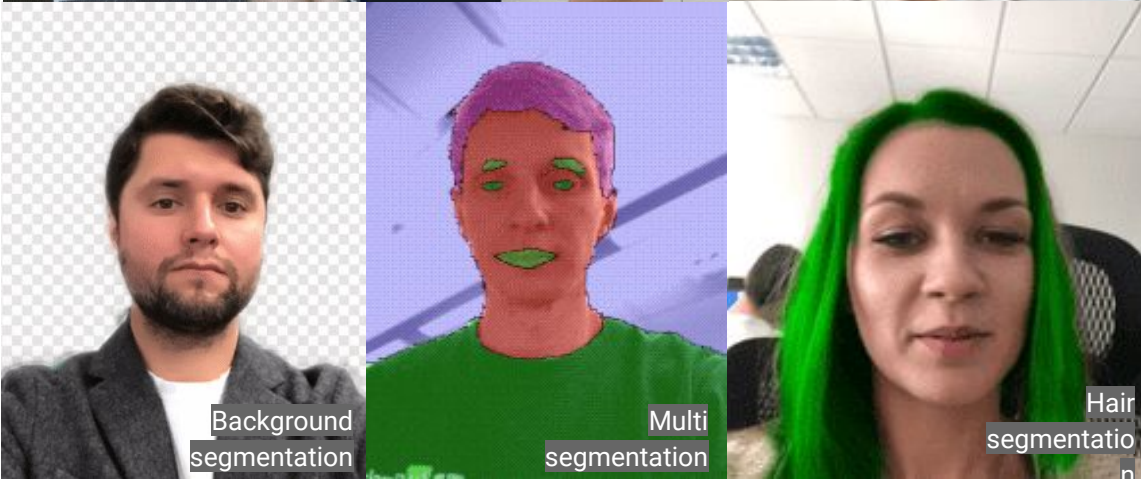
# Focus on people: Real-time on Mobile

<https://developers.google.com/ar/develop/java/augmented-faces/>



Real time:

- Hand tracking
- Face tracking
- Expression parsing



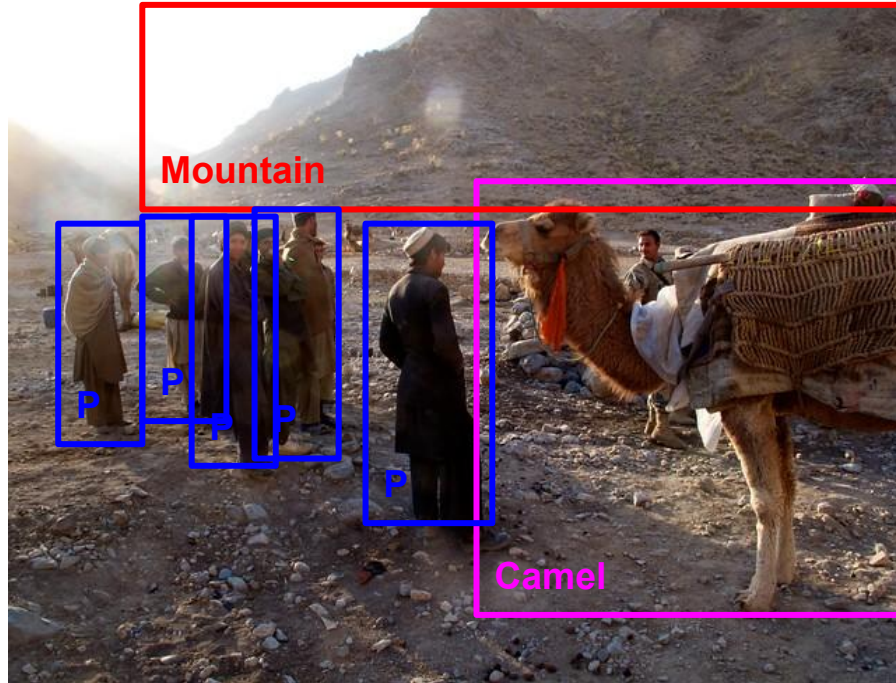


# Action Recognition

Object Recognition

!=

Action Recognition



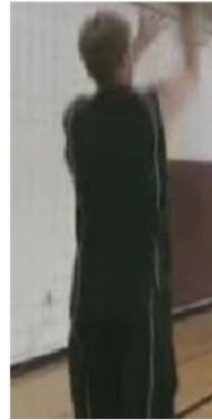
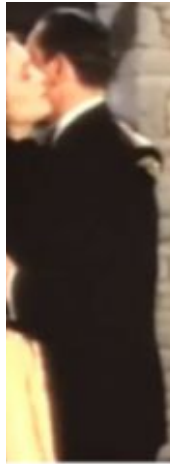
Examples of "Paint" action in AVA

[Sun, Ross, Vondrick, Pantofaru, et al.]

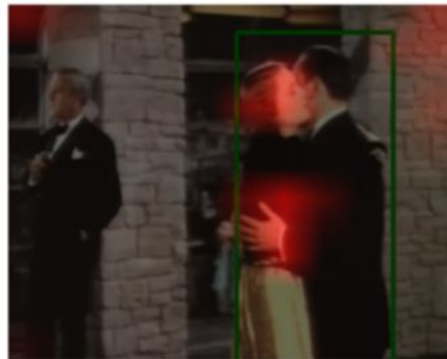
Google | YouTube

# Actor-Centric Relation Network (ACRN)

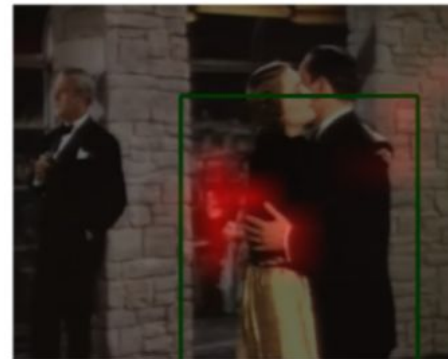
- Faster RCNN looks only at the actors (appearance, pose, etc.)
- **Opportunity: model relationship between actor and other objects/people**



# Actor-Centric Relation Network (ACRN)



Hug



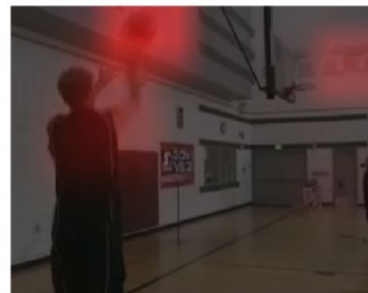
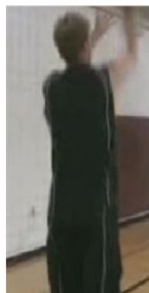
Carry



# Actor-Centric Relation Network (ACRN)

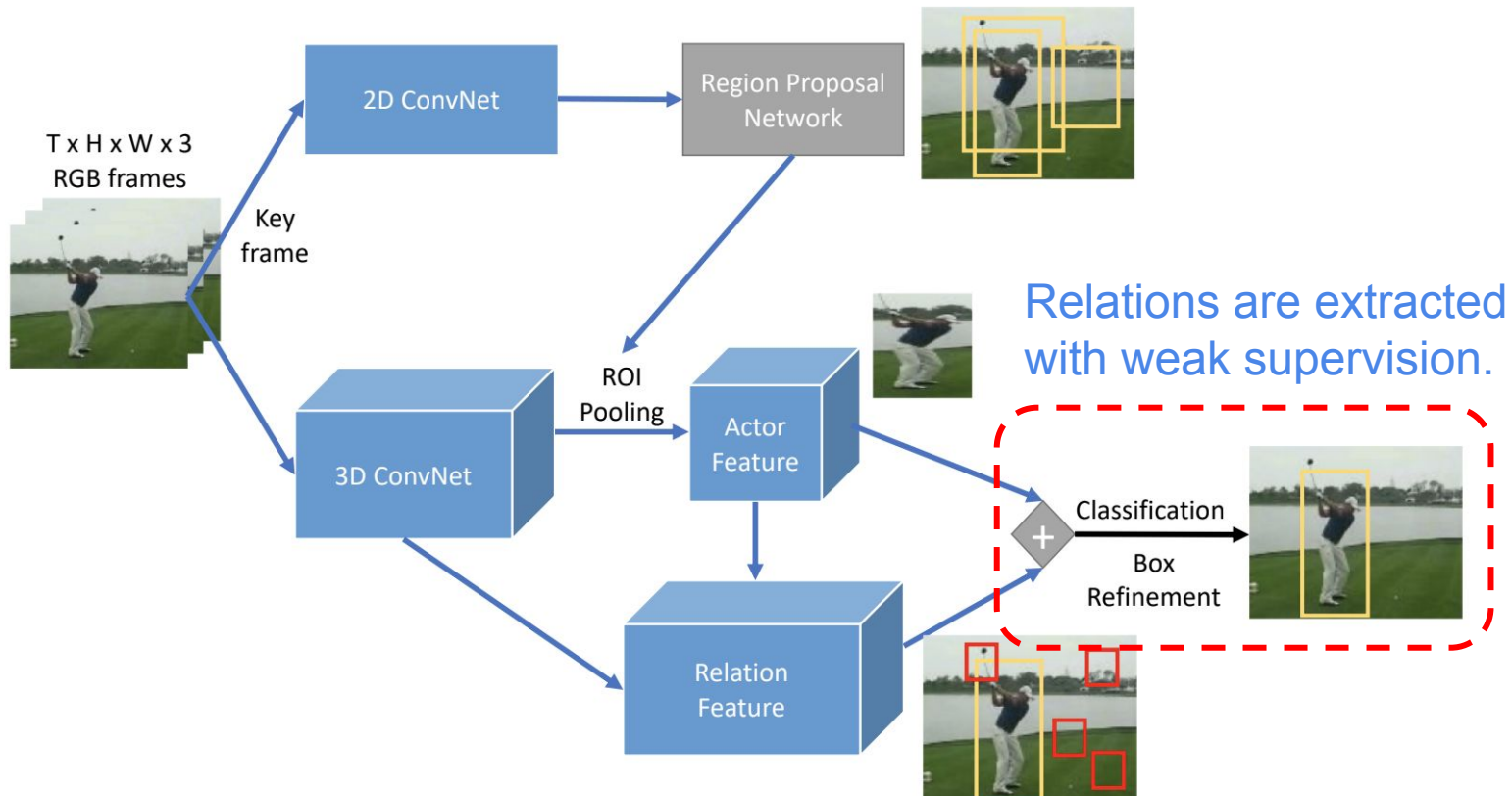


Catch ball



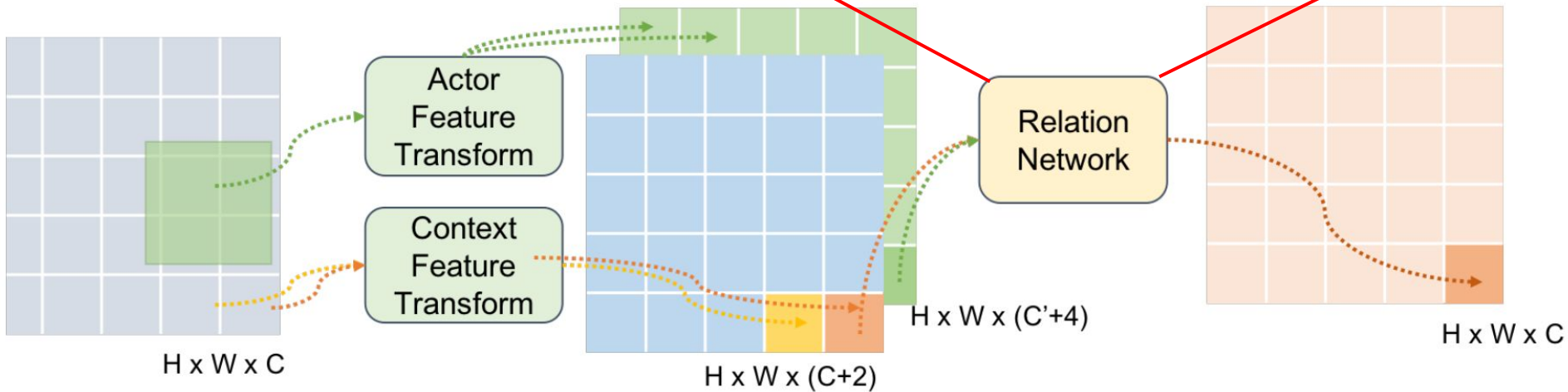
Throw ball

# ACRN Architecture



# ACRN Architecture

$$\text{ACRN}(\mathcal{A}_i, I) = f_\phi \left( \sum_{j,k} g_\theta(a_i, o_{j,k}) \right)$$



- Pairwise relation between actor and “objects”
- No explicit objectness proposals, use feature cells
- Implemented as 1x1 convolutions

# Visualizations



shoot ball



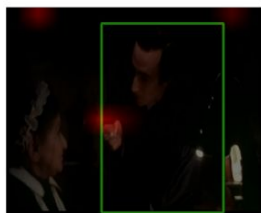
kick ball



pour



# Visualizations



smoke



listen



hug



kiss



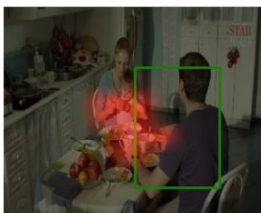
fight



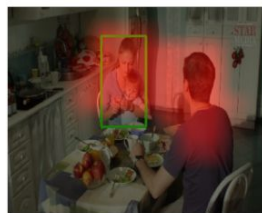
watch



eat



listen



grab



bend



read



sit





# Understanding: Starting from...

## Classification

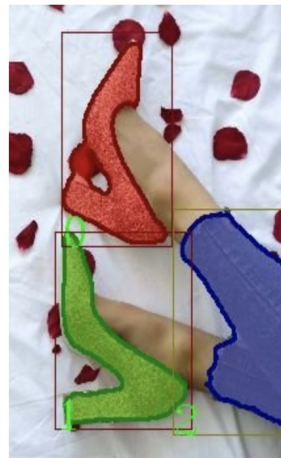
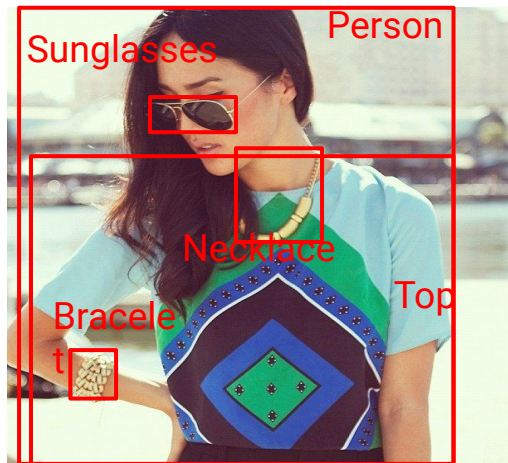
Images  $\Rightarrow$  Labels



Cake (0.93669)  
Dessert (0.91911)  
Birthday (0.89697)  
Child (0.89183)  
Fondant (0.88340)  
Birthday cake (0.87525)

## Detection

Images  $\Rightarrow$  Labelled Boxes or Regions

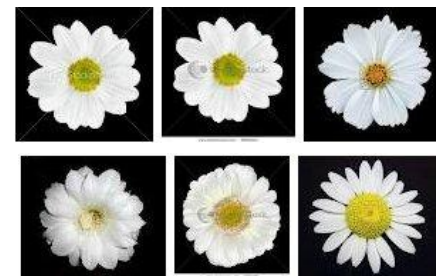


## Embedding

Images  $\Rightarrow$  Features



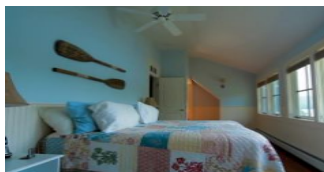
*Similar Images and Labels*



# Understanding: ... adding in...

## 3D perception

Images  $\Rightarrow$  3D relationships

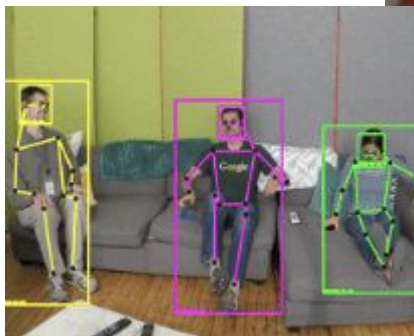
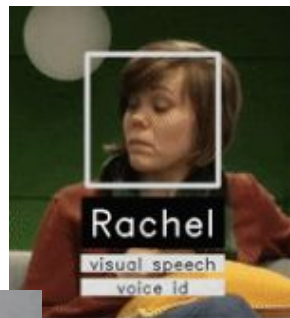
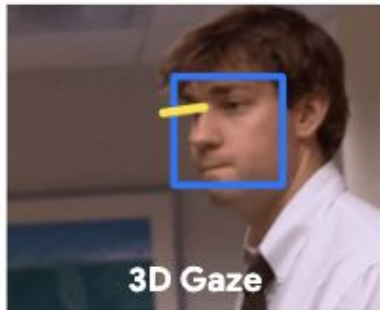


single image



## Person-centric Models

Focus, Pose, Speaker Models



## Action/Interaction Recognition



shoot ball



kick ball

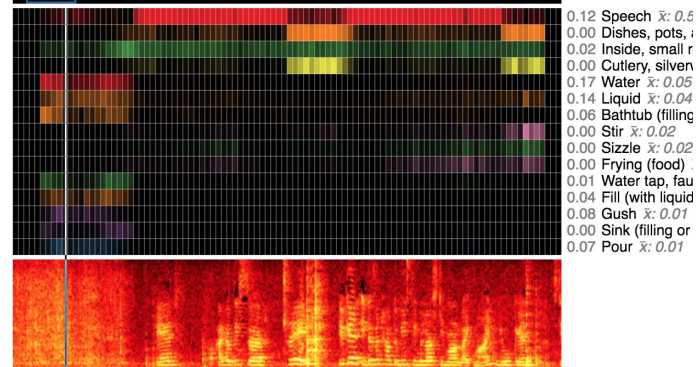


pour





# Much of video labeling/segmentation



different than understanding?  
not focused on intent (or significance)

# Text understanding: more focused on intent



Google

Thank you for filing your 1st patent application!

Your innovation is a key piece to Google's success.

The Google Patent Team

Comprehensive OCR for lots of languages



Lens Dining



Lens Tip Calculator



# Machine Translation: All about both focus and context

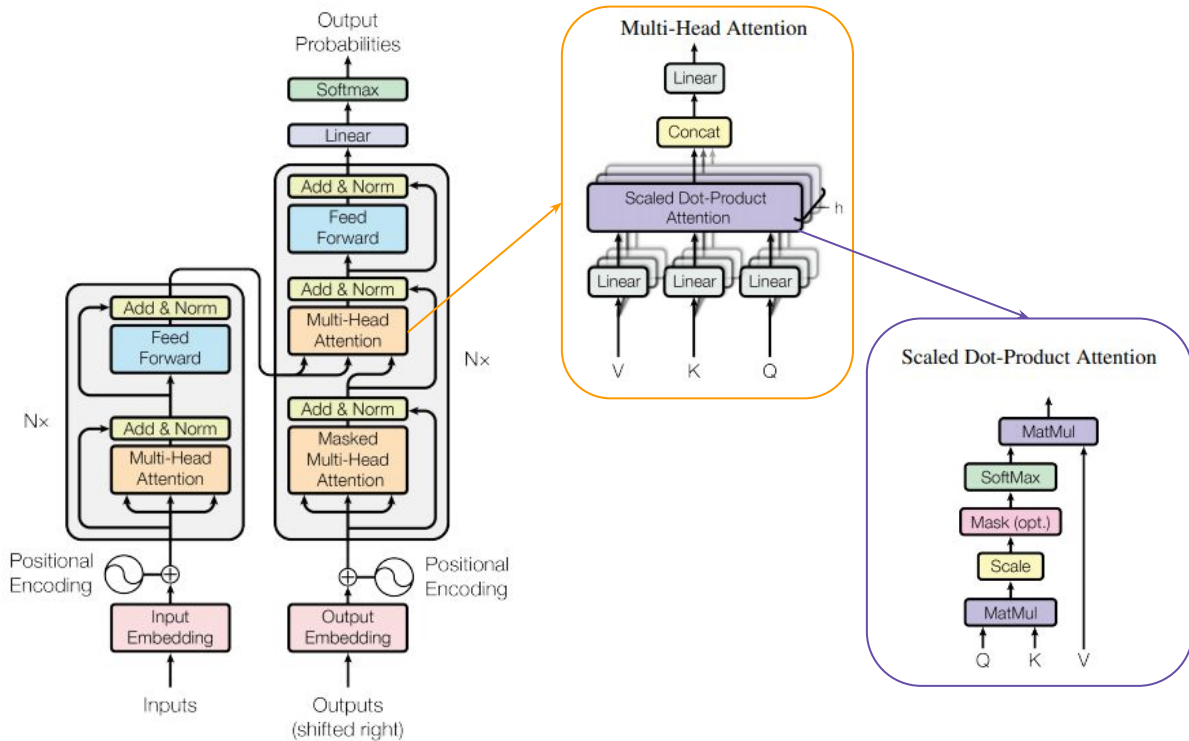
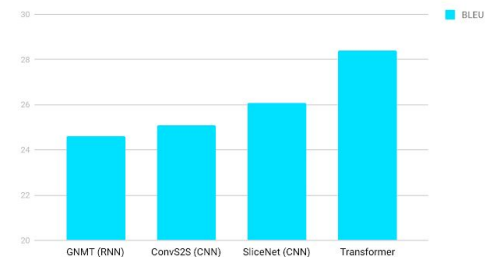


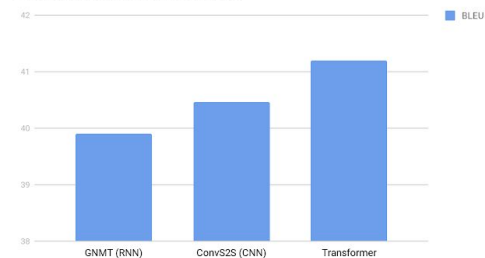
Figure 1: The Transformer - model architecture.

English German Translation quality



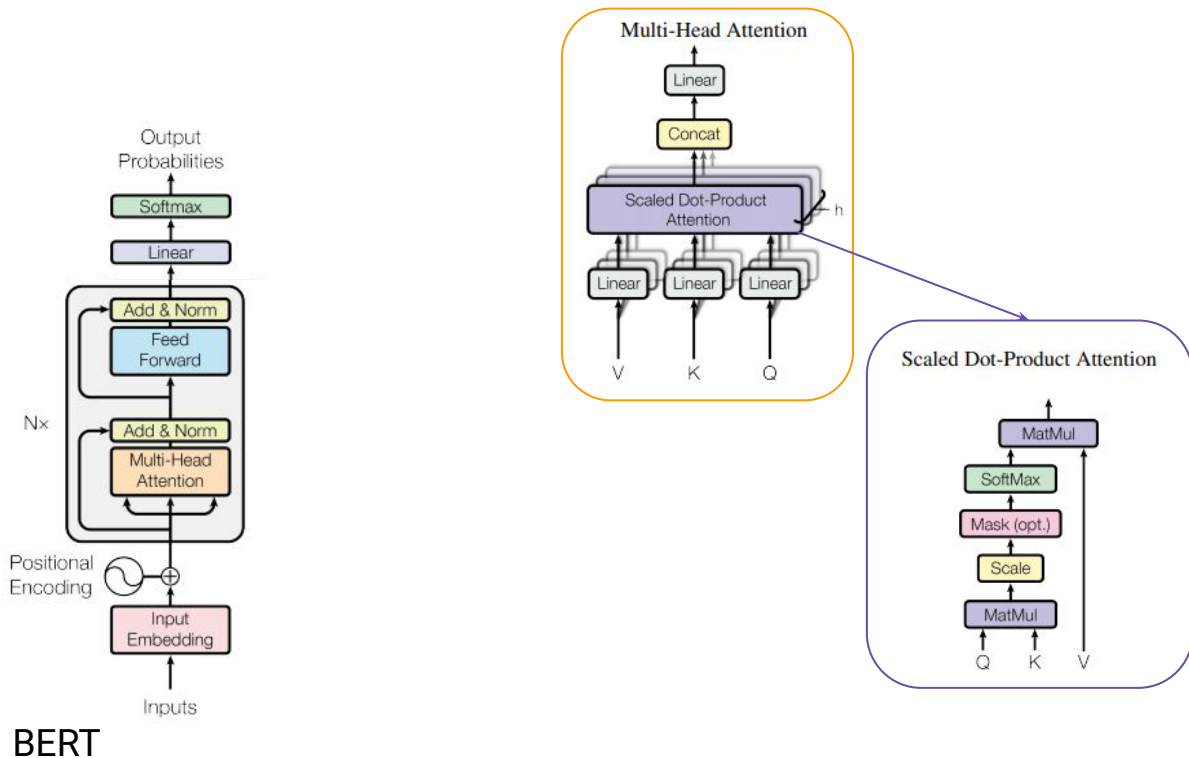
BLEU scores (higher is better) of single models on the standard WMT newstest2014 English to German translation benchmark.

English French Translation Quality



BLEU scores (higher is better) of single models on the standard WMT newstest2014 English to French translation benchmark.

# Machine Translation: All about both focus and context



## Results with BERT

To evaluate performance, we compared BERT to other state-of-the-art NLP systems. Importantly, BERT achieved all of its results with almost no task-specific changes to the neural network architecture. On **SQuAD v1.1**, BERT achieves 93.2% F1 score (a measure of accuracy), surpassing the previous state-of-the-art score of 91.6% and human-level score of 91.2%:

### SQuAD1.1 Leaderboard

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1	BERT (ensemble) Google AI Language <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	87.433	93.160
2	ninet (ensemble) Microsoft Research Asia	85.356	91.202
3	QANet (ensemble) Google Brain & CMU	84.454	90.490

BERT also improves the state-of-the-art by 7.6% absolute on the very challenging **GLUE benchmark**, a set of 9 diverse Natural Language Understanding (NLU) tasks. The amount of human-labeled training data in these tasks ranges from 2,500 examples to 400,000 examples, and BERT substantially **improves upon the state-of-the-art** accuracy on all of them:

Rank	Model	Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI	RTE
1	BERT: 24-layers, 1024 hidden, 16 heads	80.4	60.5	94.9	85.4/89.3	87.6/86.5	89.3/72.1	86.7	91.1	70.1
2	SingleTask Pretrain Transformer	72.8	45.4	91.3	75.7/82.3	82.0/80.0	88.5/70.3	82.1	88.1	56.0
3	BiLSTM+ELM+Attn	70.5	36.0	90.4	77.9/84.9	75.1/73.3	84.7/64.8	76.4	79.9	56.8

# Understanding: ... need both focus and context

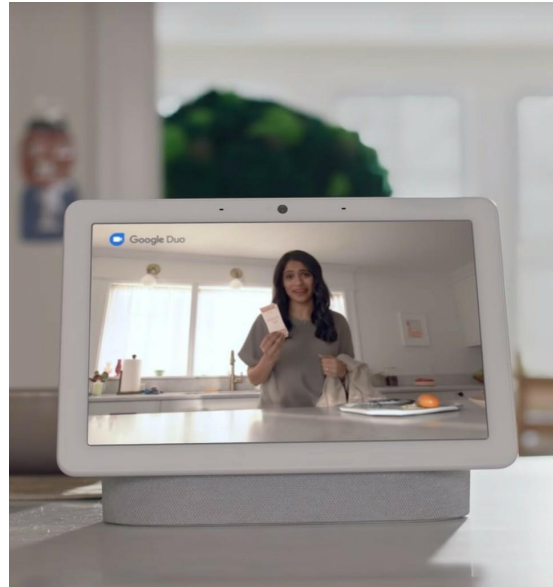
## Significance in situated video



use context



## Interaction in live video



## Intent in authored video



# Static passive monitoring cameras



- Sparse, irregular frame rate
- Power, computational, and memory constraints.
- Many images are empty
- Always looking at the same background, objects of interest often habitual



# Data Challenges

[Beery, Wu, Rathod, et al.]



(1) Illumination



(2) Blur



(3) ROI Size



(4) Occlusion



(5) Camouflage



(6) Perspective



All these images have an animal in them [Beery, Wu, Rathod, et al.]



(1) Illumination



(2) Blur



(3) ROI Size



(4) Occlusion



(5) Camouflage



(6) Perspective

# Cameras are static; Objects are habitual!

**We want per-camera models that leverage long-term temporal context to:**





# Cameras are static; Objects are habitual!

We want per-camera models that leverage long-term temporal context to:

- 1. Ignore salient false positives



These rocks have not moved in a month. *Probably not animals.*



# Cameras are static; Objects are habitual!

We want per-camera models that leverage long-term temporal context to:

- 1. Ignore salient false positives
- 2. Improve per-location object classification



Probably the same species; If we're confident about one, that should help us classify the other



# Our approach (high level)



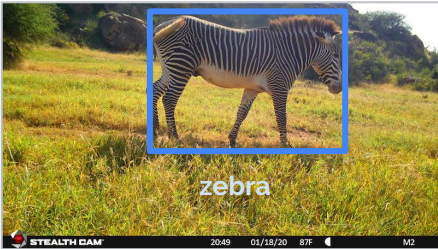
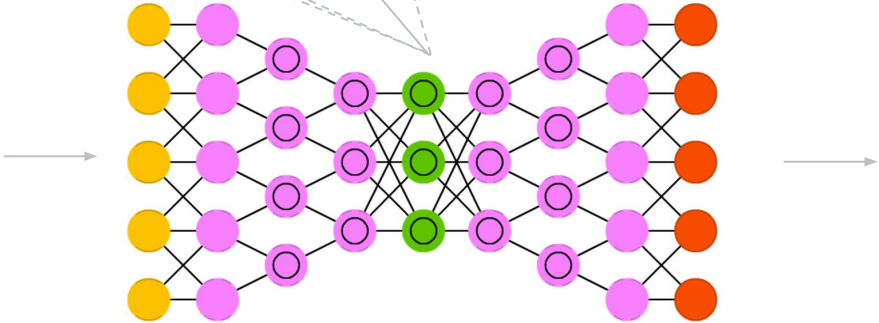
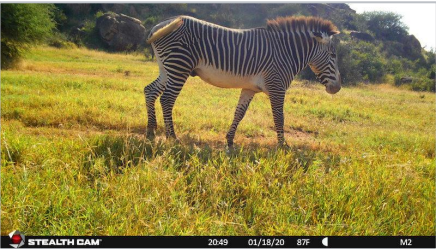
1. Build (nonparametric) per-camera model representing what a camera has seen (a.k.a. *memory bank*)



# Our approach (high level)

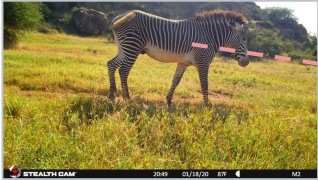
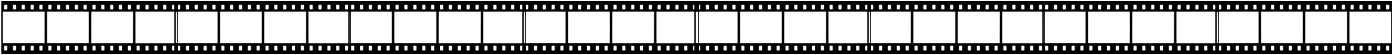


- 1. Build (nonparametric) per-camera model representing what a camera has seen (a.k.a. *memory bank*)

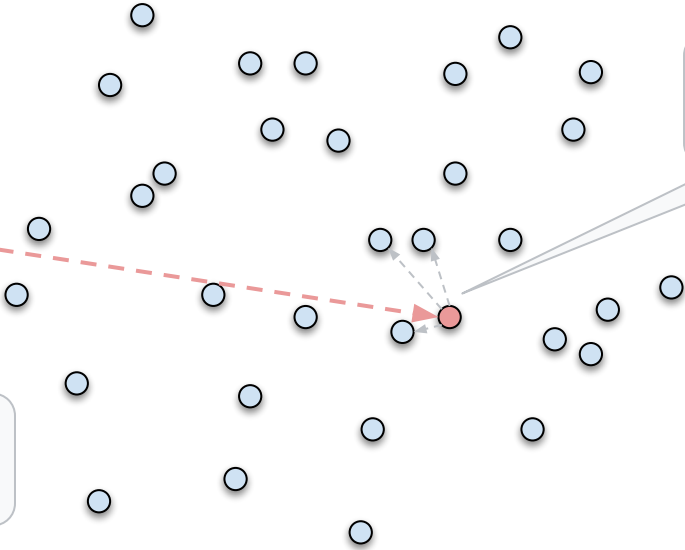


- 2. Give model running on current frame a way to reference into the memory bank

# Aggregating Features from Memory Bank: Simplified



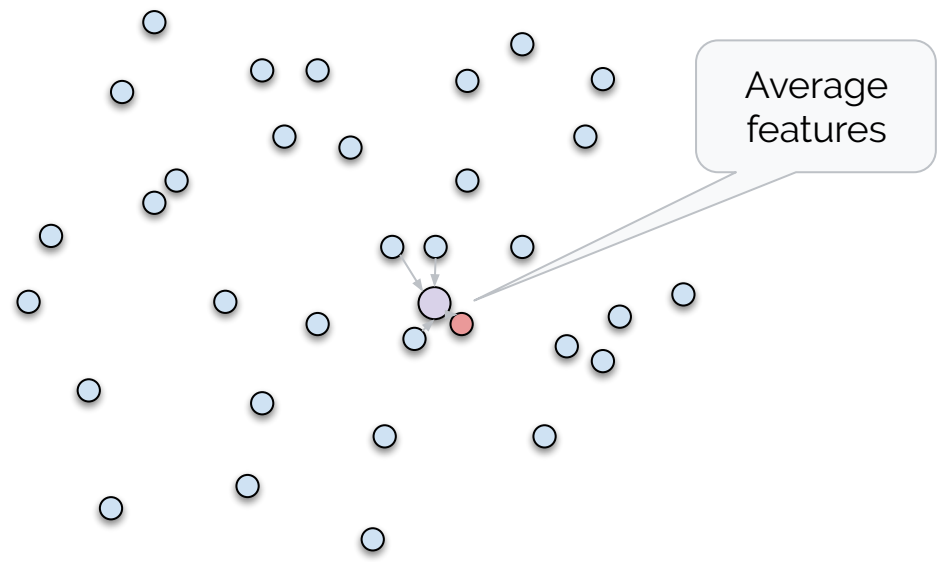
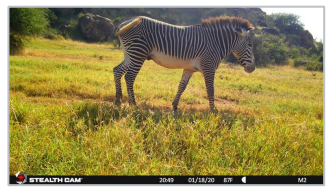
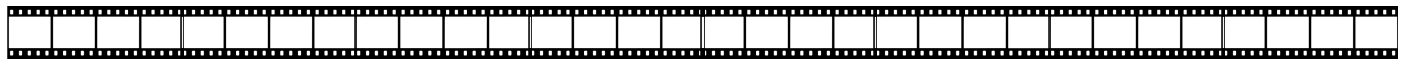
Project object proposals to embedding space



Nearest Neighbors

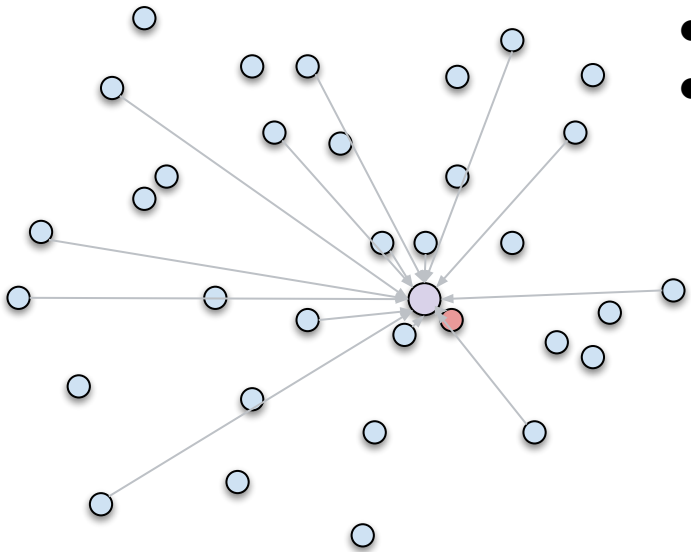
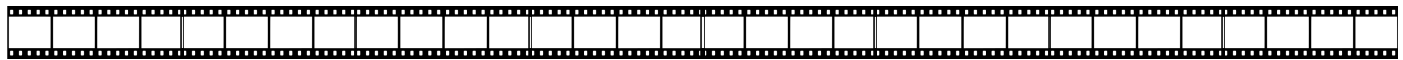
Contextual Embedding Space

# Aggregating Features from Memory Bank: Simplified

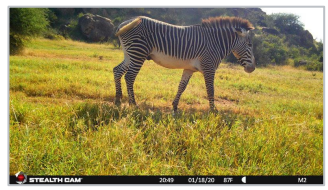


Contextual Embedding Space

# Aggregating Features via Attention



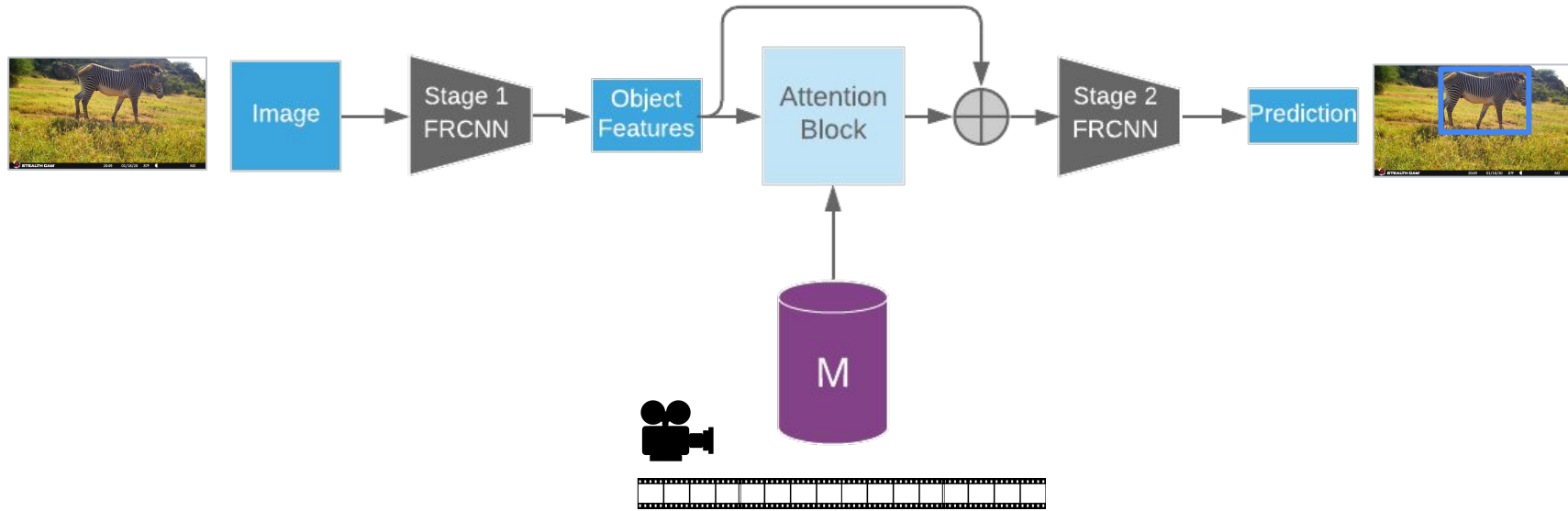
- Learn similarity metric
- Weighted average over entire memory bank



Contextual Embedding Space

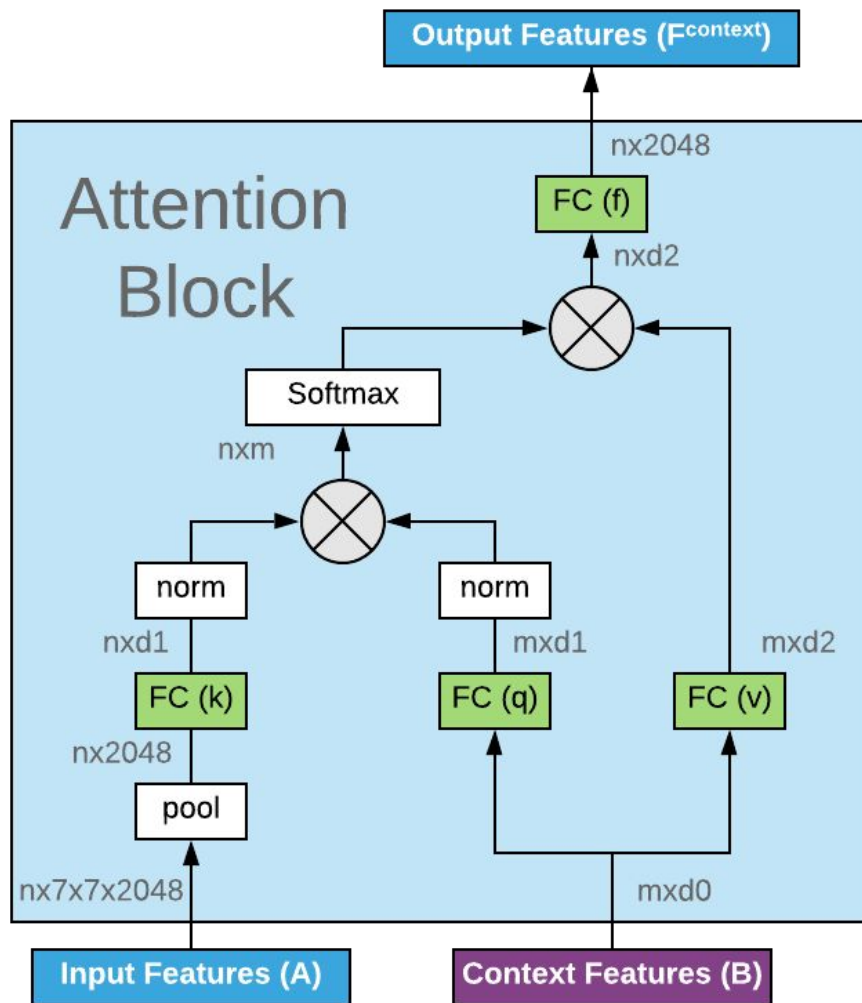


# Context R-CNN Architecture

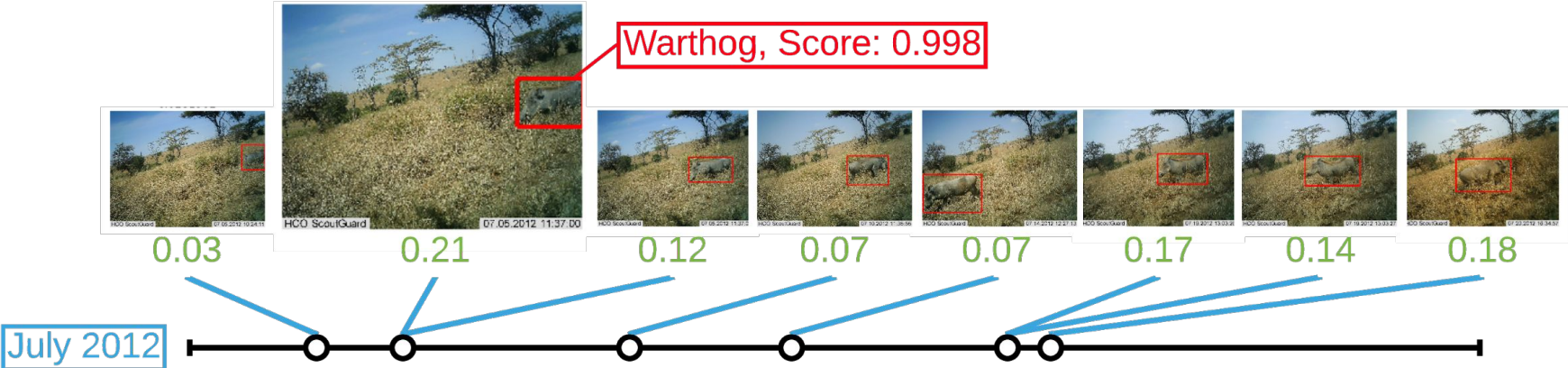


Context is incorporated based on relevance

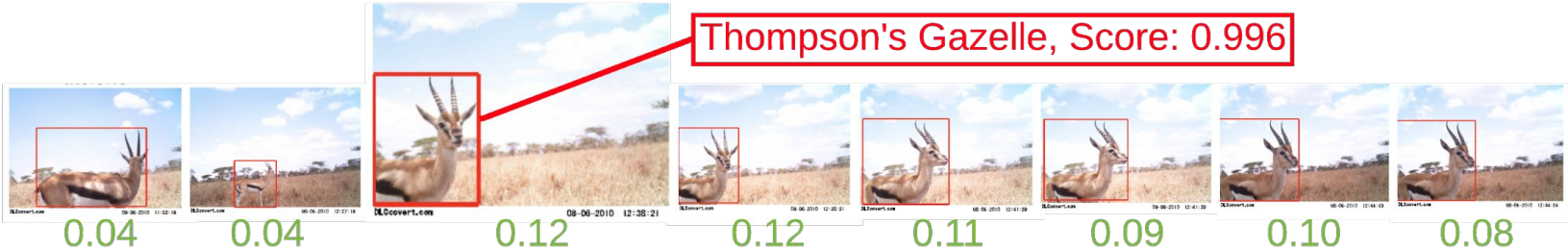
[Beery, Wu, Rathod, et al.]



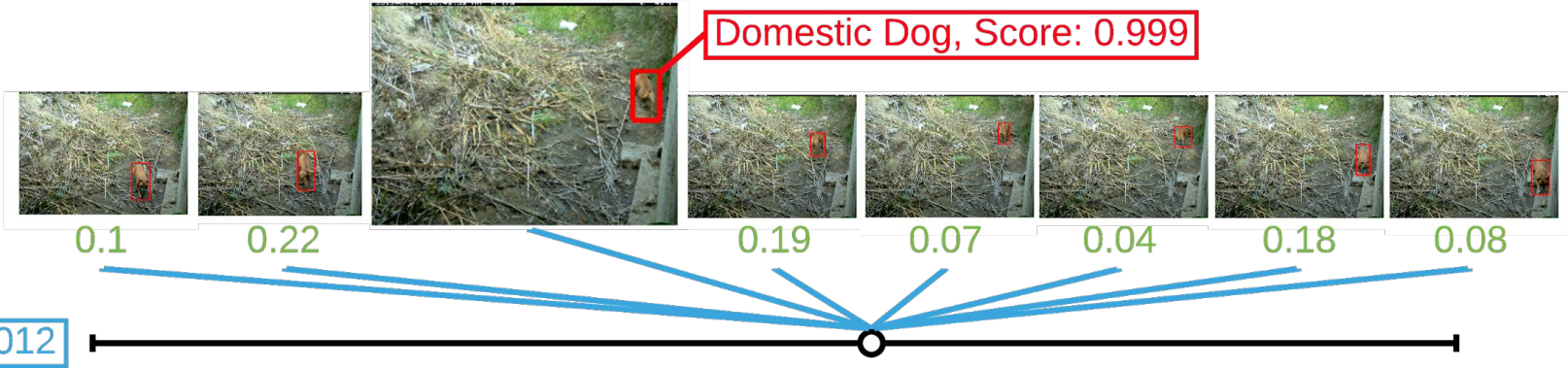
# Attention is temporally adaptive to relevance



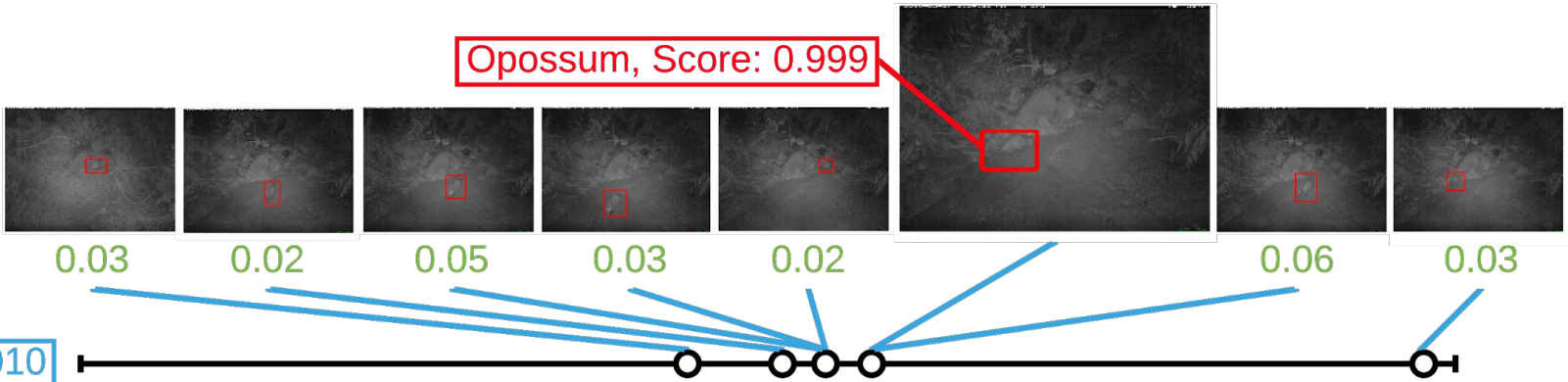
July 2012



August 2010



January 2012



May 2010



# Background classes are learned without supervision





# Authored video content: understand intent!

<https://www.blog.google/outreach-initiatives/accessibility/get-image-descriptions/>

## Machine-generated image descriptions



Machine-generated description for this image: "Appears to be: Person playing guitar on the sofa."

Person playing guitar on the sofa



Machine-generated description for this image: "Appears to be: Fruits and vegetables at the market."

Fruits and vegetables at the market



# Early Steps:

## Exploiting Speech to Train Video Representations



- Cross-modal weak supervision: ASR  $\Leftrightarrow$  vision  
(exploits co-occurring but noisy speech to supervise representation)
- Multimodal input: audio + video
- Some progress towards longer video sequences



# Video BERT



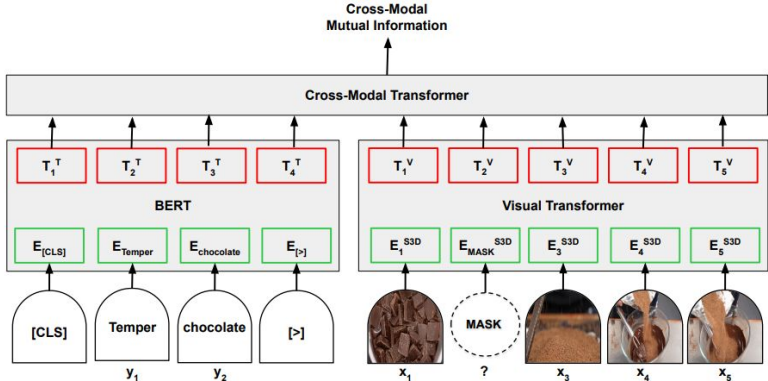
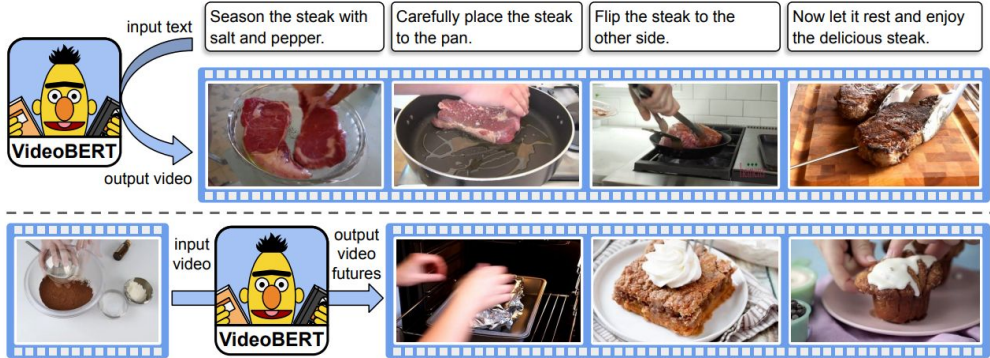
- Example 1: Given recipe text, generate sequence of visual tokens (retrieved from different videos)

# Video BERT



- Example 2: Given a video token, predict possible futures (cocoa & flour mix can get baked and turn into a brownie or cupcake)

# Joint video and language representations



- Understand long videos with **self- (time) and cross-modal (ASR) supervisions**.
- Leverage powerful models (BERT) and pre-training tasks (masked LM).
- Cross-modal applications: zero-shot action classification, action anticipation, etc.
- Opportunities:
  - Automatic video data mining given large vocabulary (Video Search timed anchors).
  - Generic feature vectors for long videos (VCA).

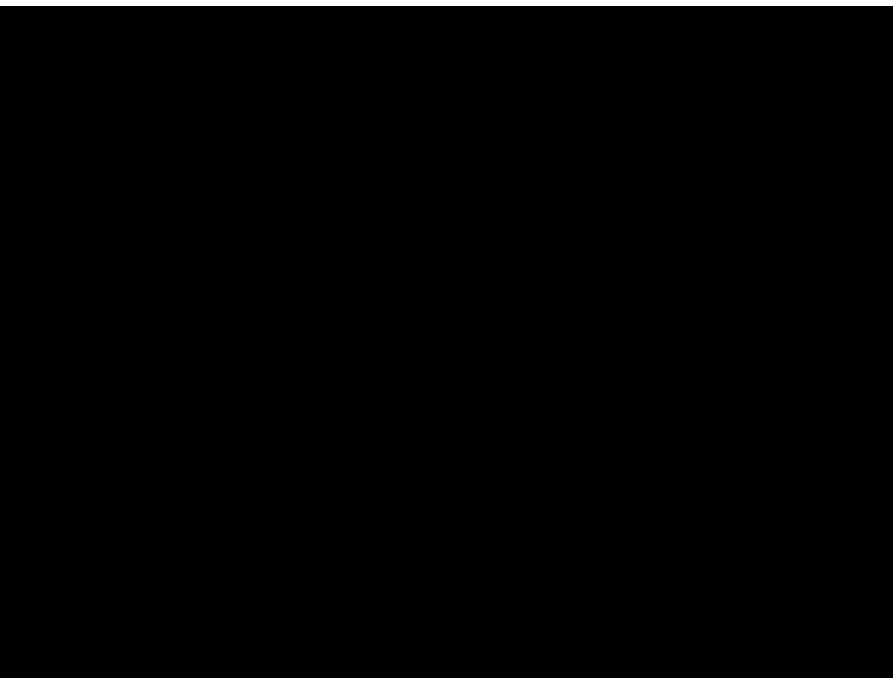
# Live interactions



Significance  
and intent both  
come into play

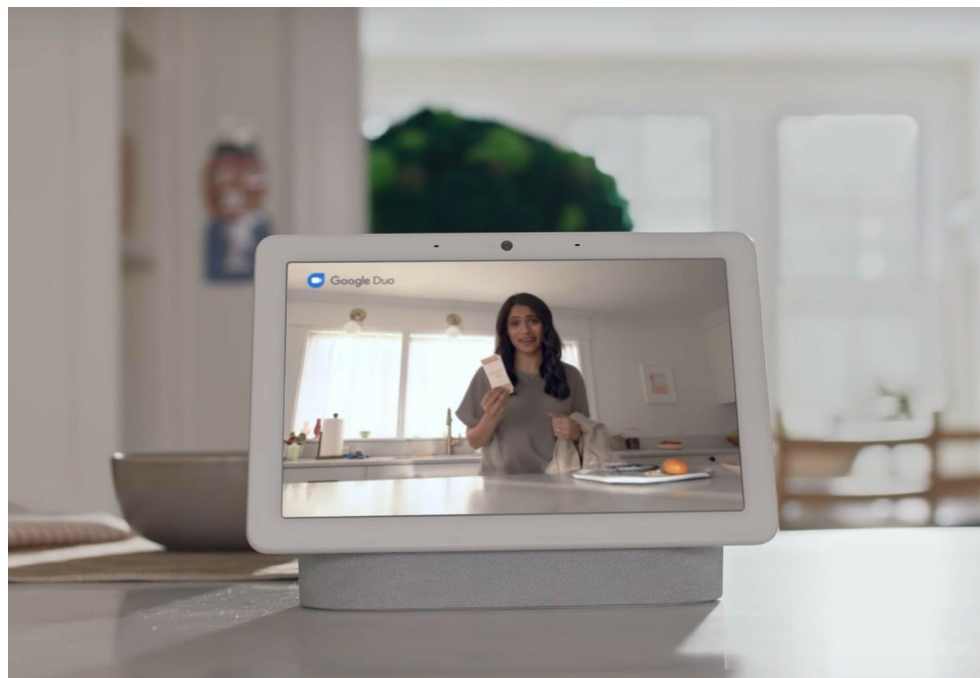


# Early Steps in Situated Perception



BodyPix - Real-time Segmentation in your Browser

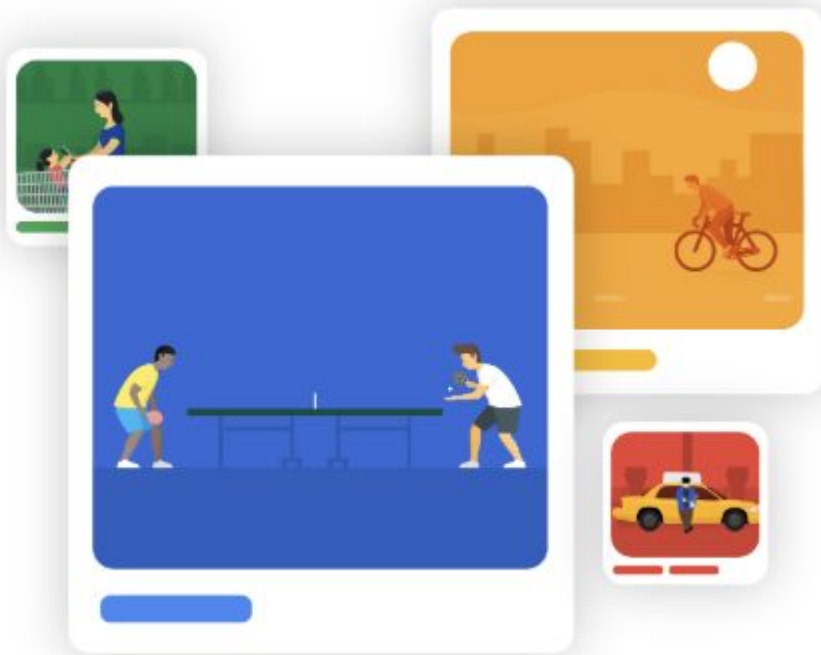
<https://github.com/tensorflow/tfjs-models/tree/master/body-pix>



Natural gestures & smart camera in Google Nest Hub Max

Image credit: <https://www.dailydot.com/wp-content/uploads/2019/05>

# Need to understand what to leave out



Media retargeting  
as a (small) window  
into this problem

# Retargeting in time



Video



Preview (6s)



Summary (9s)



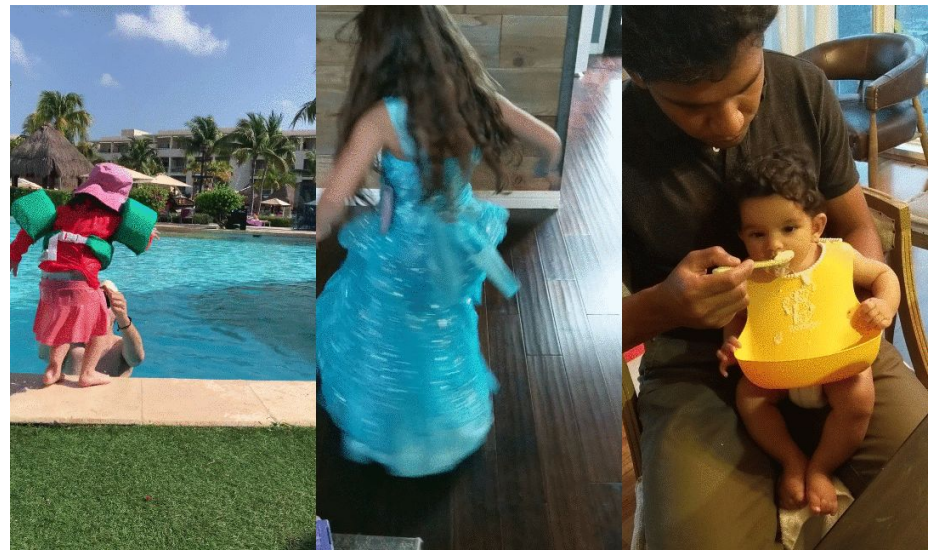
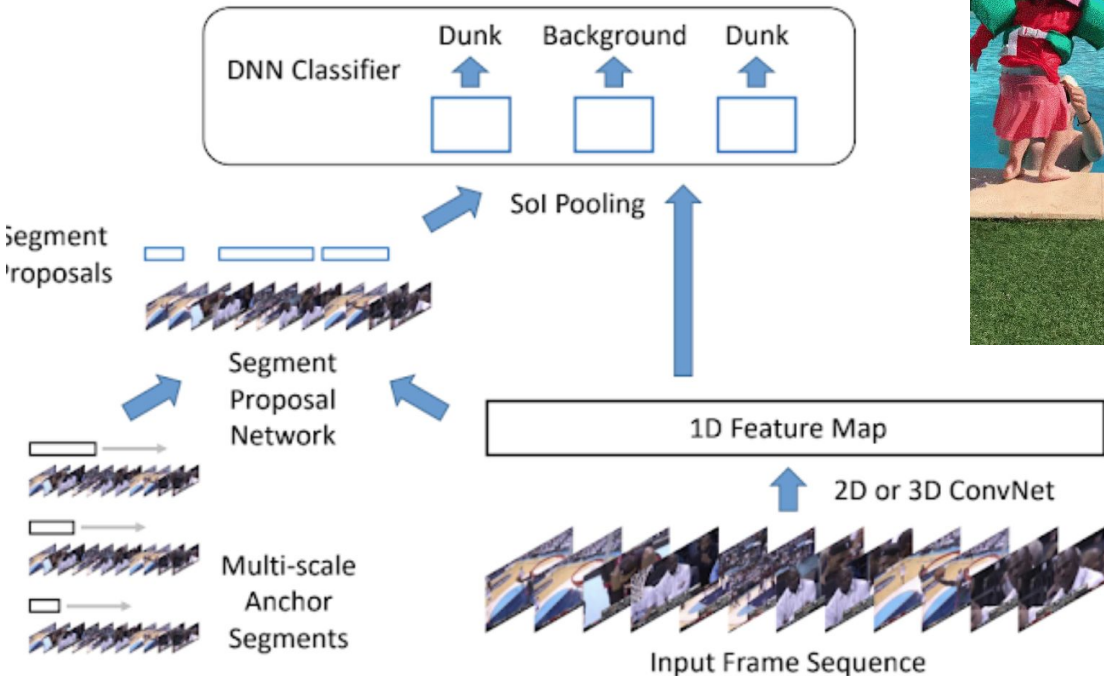
Video ad (1 minute)



Bumper ad

# Retargeting in time

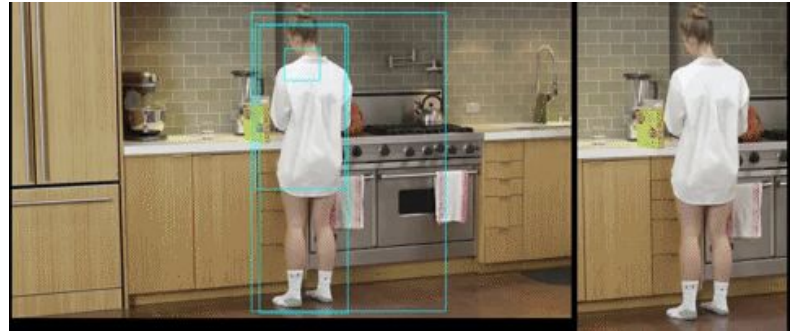
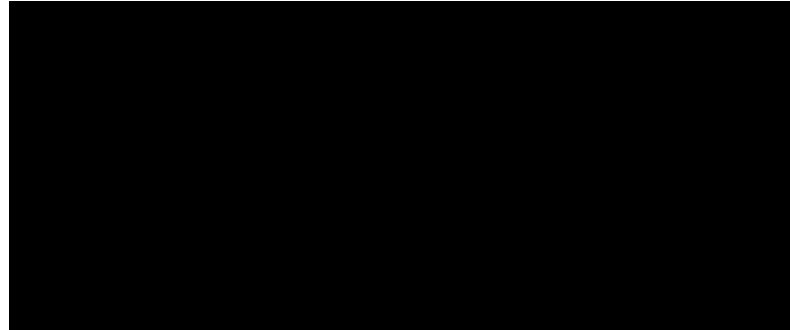
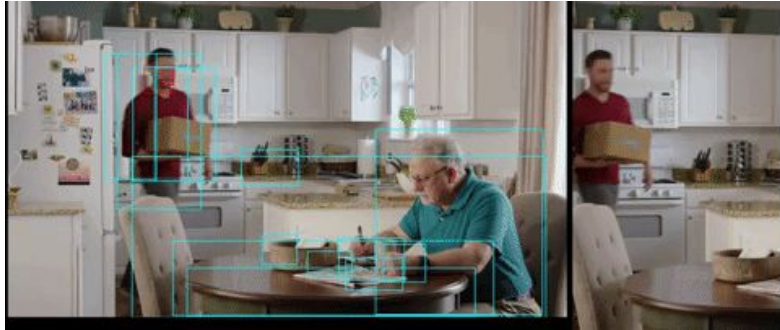
[Chao, Vijayanarasimhan, Seybold, et al.]



Special moments in video



# Retargeting in space



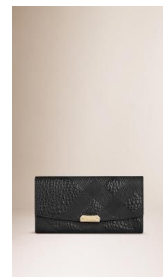
L: original landscape; R: AutoFlip

<https://www.blog.google/products/ads/level-your-gaming-business-new-innovations-apps>

# Retargeting in Space

Deep-learned, aesthetic cropping

Improves quality and generates fewer bad crops by training on millions of professional-quality photos



# Directions in ML

Understanding **intent** and **importance** in audio and video

- detection and categorization are only part way
  - also have: 3D understanding, people-centric information, action/interaction recognition
- missing piece: need to determine **what to leave out**
  - want a **synopsis** for **authored** (and situated) media
  - want less constrained interactions for **live** situations

Generating new **creative** content can help highlight shortcomings  
(as well as providing useful content)