# Deep Context for Fine-Grained Object Recognition in Crowded Images

ERAN GOLDMAN

Joint work with Prof. Jacob Goldberger,
Faculty of Engineering, BIU
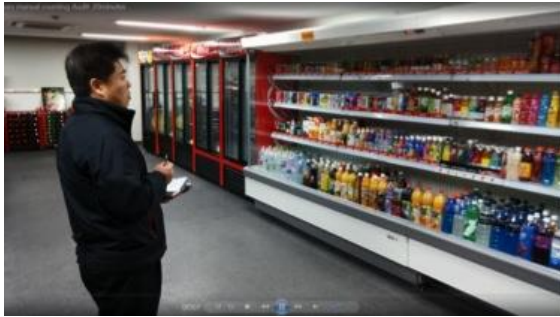
Trax
image recognition

Trax
image recognition

# AGENDA

❯ **Trax Visual Challenges**

❯ **Deep Context Embedding Architecture**

❯ **Implementation Details**

❯ **The Detection Challenge**

❯ **Summary**

**Trax** image recognition

# Traxs' Business Application



| Manual Audit | Trax Automatic Recognition Audit | |



**Slow and Expensive Inconsistent untraceable** → **Fast and Cheap Consistent Traceable** → **'Big Data' for retail**

**AVAILABILITY** **SHARE OF SHELF** **PRICING** **PROMOTIONAL ACTIVATIONS** **COMPETITVE INSIGHTS** **PLANOGRAM COMPLIANCE** **SHELF STANDARDS**

**Trax** image recognition

# Trax unlocks 'Big Data' for the retail industry

## Scale of the Data

**Market**

Market share, trade channels, segments, competition

**Retailer**

Brand champions, generic brands, range review

**Store**

Product location, assortments, shelf share

**Shelf**

NPD, POSM, Pricing, Promotions

## Scale of coverage

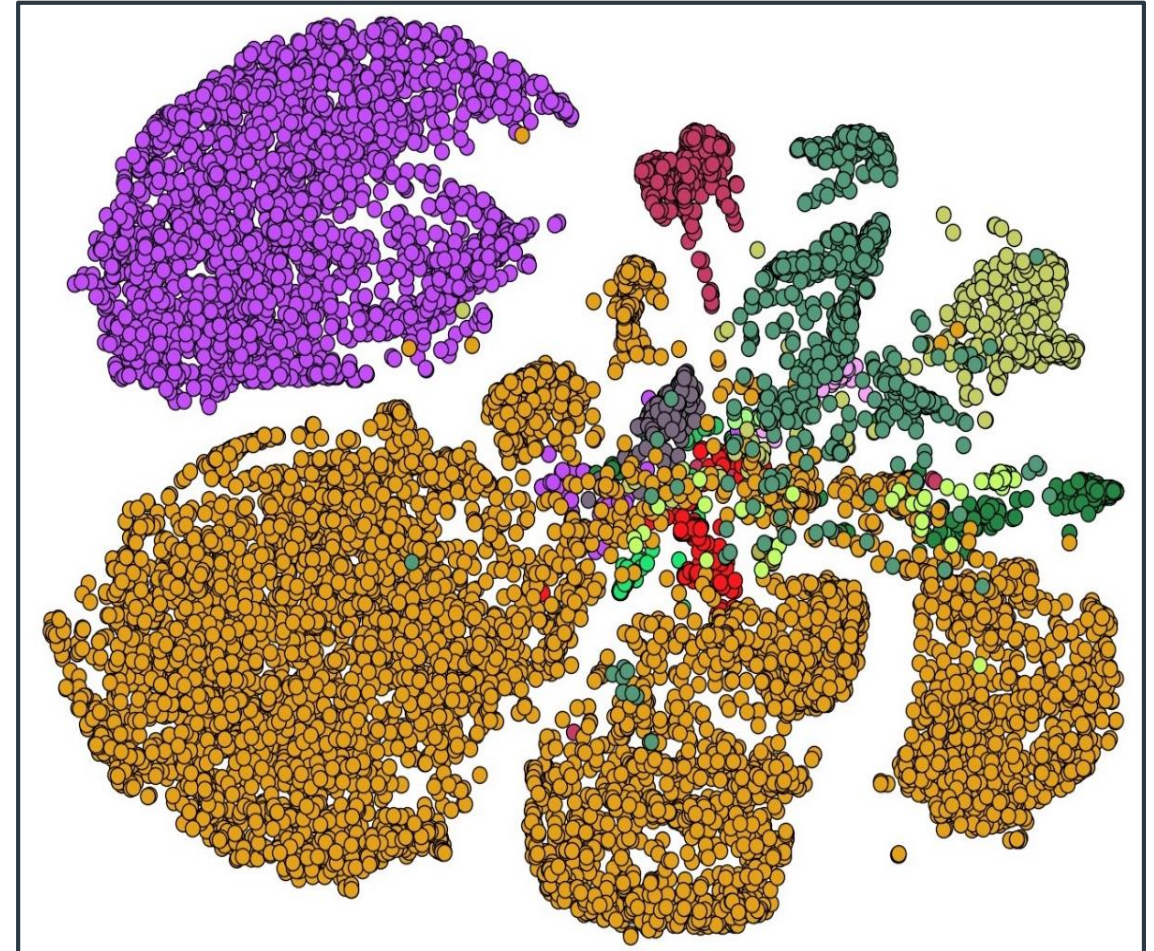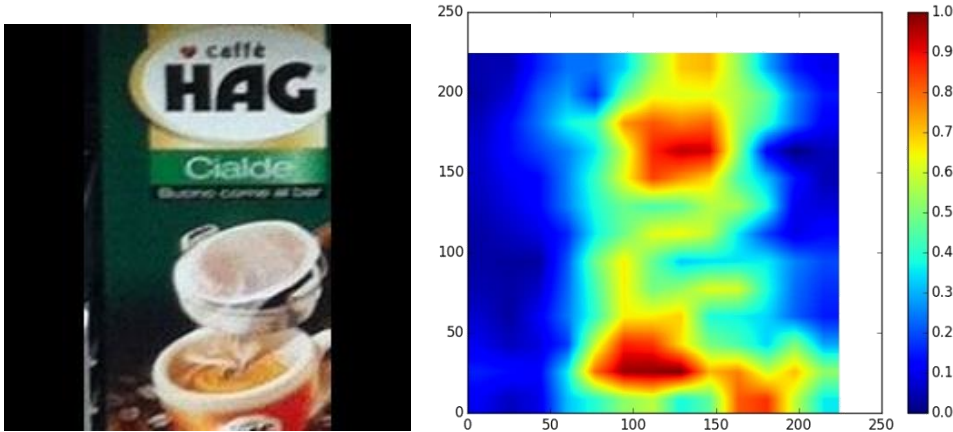# Welcome to Trax Universe

# Trax's Visual Challenges



# Classes

Fine-Grained Classification

Crowded Scene
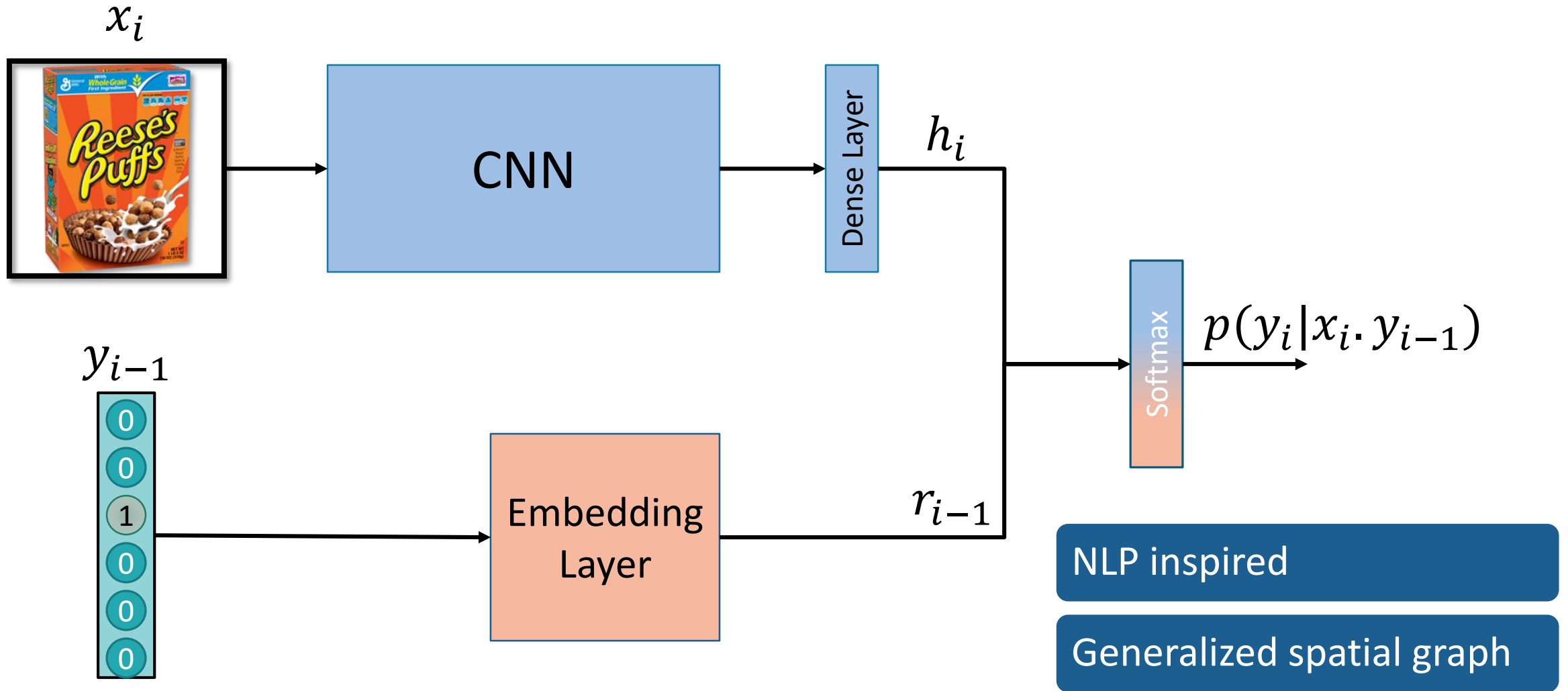
Dynamic Dataset

# Last Year- CNN: Opening The Black Box

# *AGENDA*

❯ **Trax Visual Challenges**

❯ **Deep Context Embedding Architecture**

❯ **Implementation Details**

❯ **The Detection Challenge**

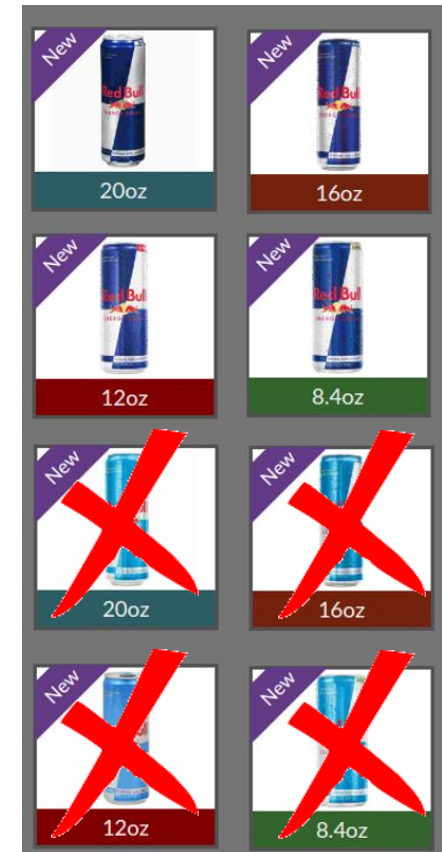❯ **Summary**

# New DNN architecture

# Which one is it?

classes

Object

**Trax** image recognition

# Which one is it?

classes

Object

**Trax**
image recognition

# Which one is it?



Object

classes

Trax
image recognition

# Which one is it?



classes

Object

Trax image recognition

# Which one is it?

classes

Flavors with red labels:

| Classic | Decaf | Lemon | Vanilla | Ginger |
|---------|-------|-------|---------|--------|

X Sizes:

| 0.35 | 0.47 | 0.5 | 0.6 | 1 | 1.25 | 2 |
|------|------|-----|-----|---|------|---|

Trax
image recognition

# Which one is it?



Object

classes

**Flavors with red labels:**

| Classic | Decaf | Lemon | Vanilla | Ginger |

Sizes:

| 0.35 | 0.47 | 0.5 | 0.6 | 1 | 1.25 | 2 |

Trax
image recognition

# Which one is it?

**Object**

**classes**

| 0.35 ltr (12 OZ) | 0.47 ltr (16 OZ) | 0.5 ltr | 0.6 ltr (20 OZ) | 1 ltr | 1.25 ltr | 2 ltr |

# Which one is it?



Object

classes

| 0.35 ltr (12 OZ) | 0.47 ltr (16 OZ) | 0.5 ltr | 0.6 ltr (20 OZ) | 1 ltr | 1.25 ltr | 2 ltr |

Trax
image recognition

# Which one is it?

Object

classes



| 0.35 ltr (12 OZ) | 0.47 ltr (16 OZ) | 0.5 ltr | 0.6 ltr (20 OZ) | 1 ltr | 1.25 ltr | 2 ltr |

Trax
image recognition

0.35 ltr (12 OZ)

0.47 ltr (16 OZ)

0.5 ltr

0.6 ltr (20 OZ)

1 ltr

1.25 ltr

2 ltr

Trax
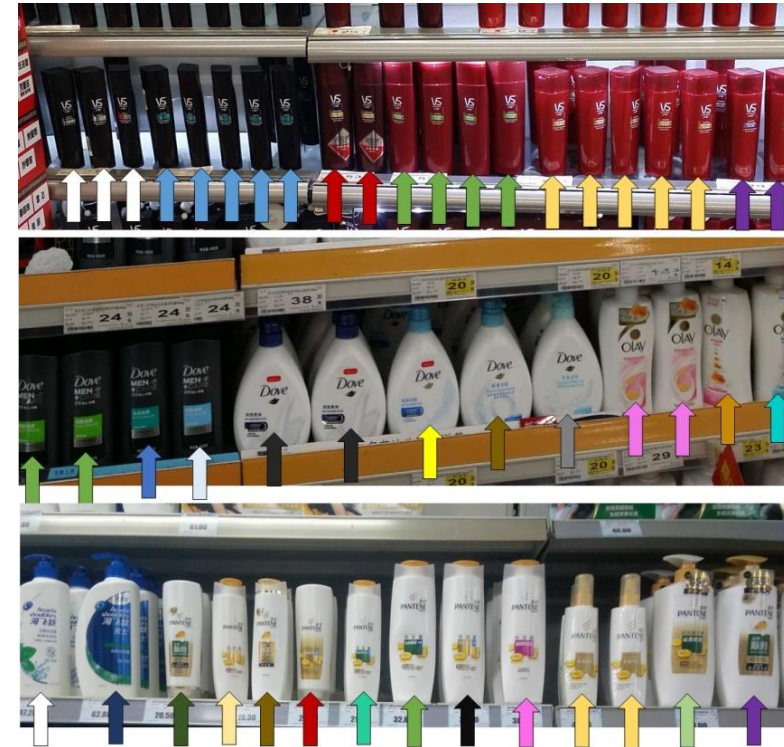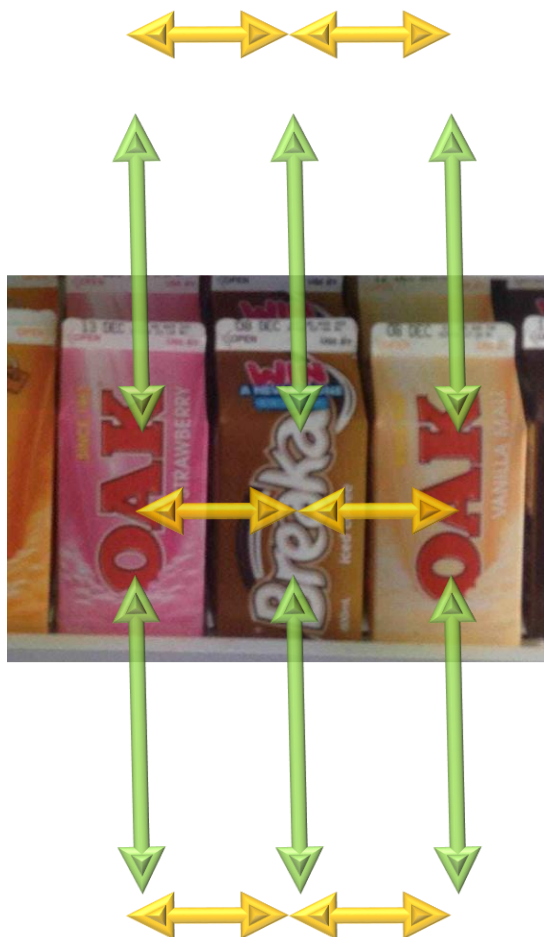image recognition

# Natural Language Processing Analogy

SENTENCES

SHELVES

A tall glass of orange juice is the very image of refreshment.

Overwatering an orange tree can cause the leaves to turn pale.

The fiber content in an orange peel makes you feel fuller after eating.

**Trax** image recognition

Trax
image recognition

# AGENDA

❯ **Trax Visual Challenges**

❯ **Deep Context Embedding Architecture**

❯ **Implementation Details**

❯ **The Detection Challenge**

❯ **Summary**

**Trax** image recognition

# Extended "Sentence"

## Model image as a graph



Detection Node

Left-Right Relation Edge

Top-Bottom Relation edge

**Trax** image recognition

# Classification



$$P(y_8|x_8)$$

# Classification



$$P(y_8 | x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8)$$

Trax
image recognition

Joint Probability Distribution

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$ $x_8$

$$P(Y/X) \propto \prod_{t=1}^{n} \varphi(y_t | x_t, y_{t-1}) \propto \prod_{t=1}^{n} e^{[r_{y_{t-1}}^T, h_t^T]W_{y_t} + b_{y_t}}$$

Trax
image recognition

Marginal Distributions

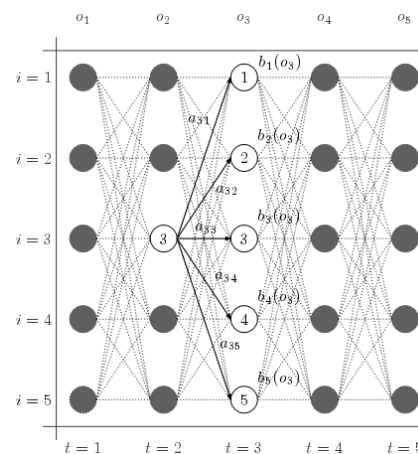$$P(y_t = i|X) = \sum_{Y/y_t} P(Y/X) =$$

$$\sum_{Y/y_t} P(y_1, y_2, \cdots, y_{t-1}, i, y_{t+1}, \cdots, y_n | x_1, x_2, \cdots, x_n) \propto$$

$$\sum_{Y/y_t} \prod_{\tau=1}^{n} e^{h_\tau^T u_{y_\tau} + b_{y_\tau}} e^{r_{y_{\tau-1}}^T q_{y_\tau}} = \alpha(y_t) \cdot \beta(y_t)$$
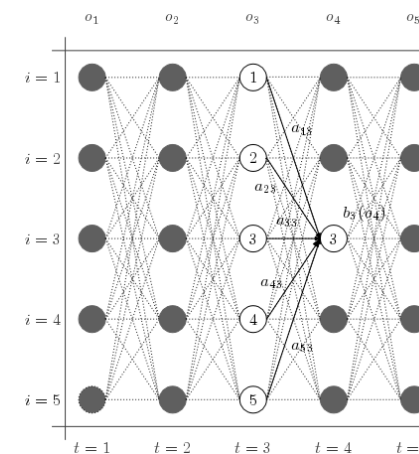
Trax
image recognition

Dynamic programming:
Forward–backward algorithm
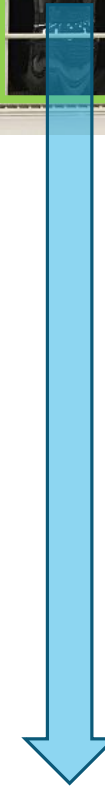
$$P(y_t|X) = \frac{1}{Z}\alpha(y_t)\beta(y_t)$$
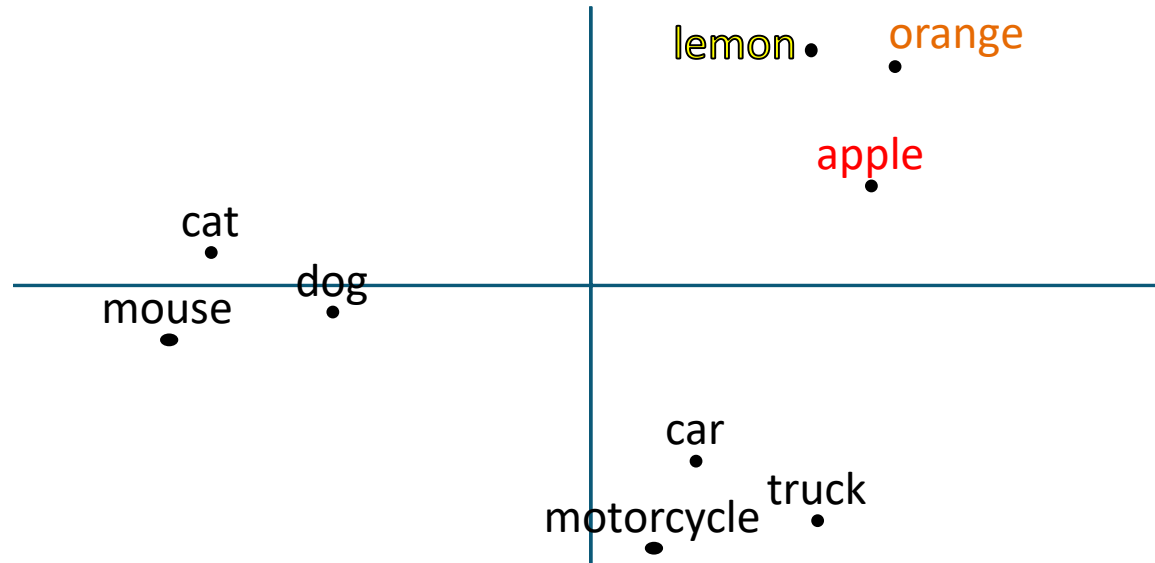
$\alpha(y_t)$

$\beta(y_t)$

Classification: $P(y_i|x_i)$

Spanning Tree Belief-Propagation

Classification: $P(y_i|X)$

Trax
image recognition

# Word Embedding $2^{nd}$ Order Similarity

*Orange* juice     *Lemon* juice

*Orange* water     *Lemon* water

*Orange* tree     *Lemon* tree

*Orange* seed     *Lemon* seed

*Orange* peel     *Lemon* peel

*Orange* pulp     *Lemon* pulp

*Orange* orchard     *Lemon* orchard

*Orange* milkshake     *Lemon* milkshake

*Orange* pie     *Lemon* pie

*Orange* soda     *Lemon* soda

*Orange* plantation     *Lemon* plantation

*Orange* drink     *Lemon* drink

*Orange* flavored     *Lemon* flavored

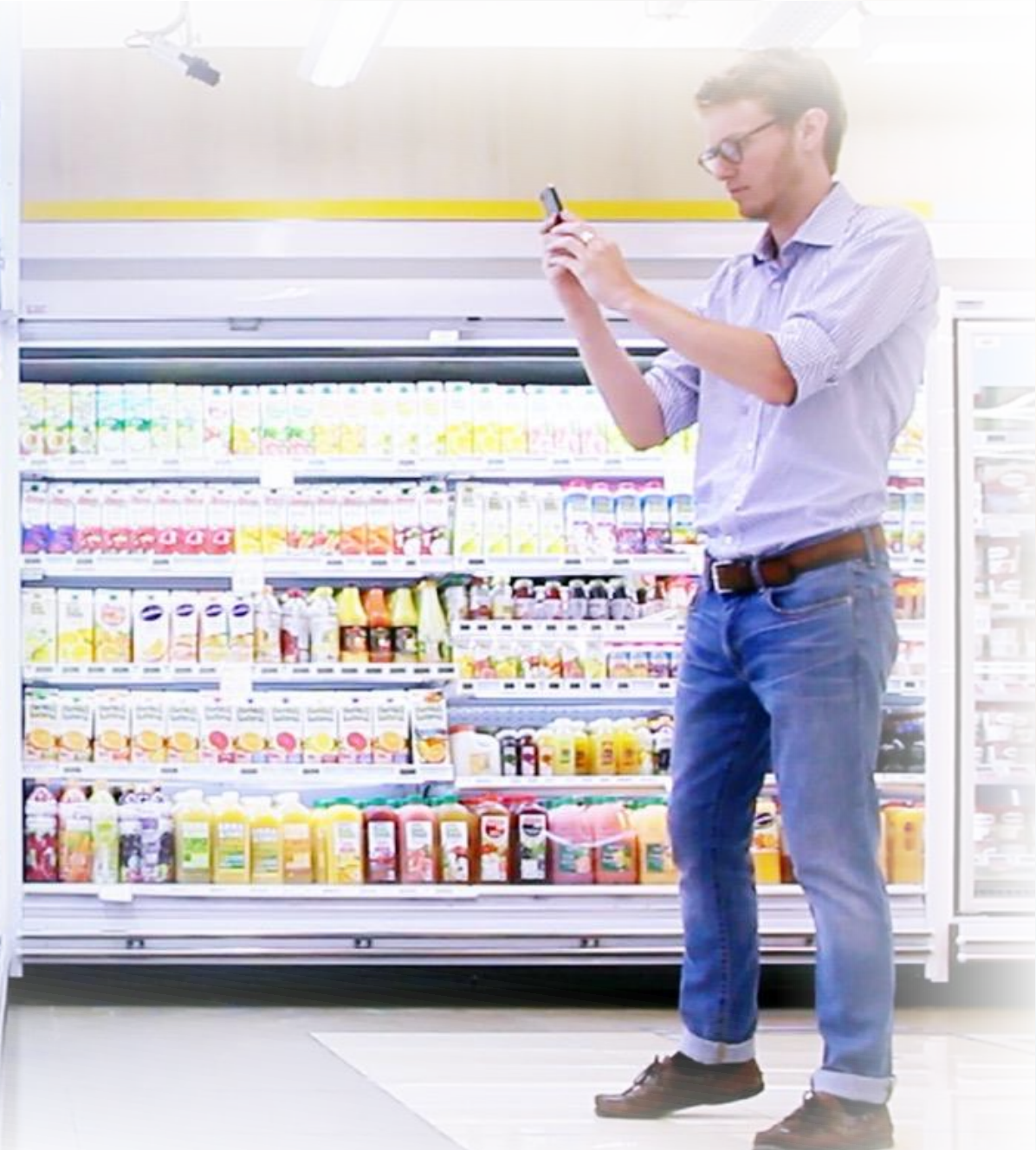*Orange* cider     *Lemon* cider

*Orange* vinegar     *Lemon* vinegar

Trax
image recognition

# Product Embedding $2^{nd}$ Order Similarity

# *AGENDA*

❯ **Trax Visual Challenges**

❯ **Deep Context Embedding Architecture**

❯ **Implementation Details**

❯ **The Detection Challenge**

❯ **Summary**

# Correct Detection

# Missing an object

## you shall ~~not~~ pass

**Trax** image recognition

# Splitting an object

## ....the fireworks on the fourth of July

**Trax** image recognition

# State of The Art Detectors

"YOLO imposes strong spatial constraints on bounding box predictions"

*Redmon, Joseph, et al. "You only look once: Unified, real-time object detection."*

"SSD is very sensitive to the bounding box size"

*Liu, Wei, et al. "SSD: Single shot multibox detector."*

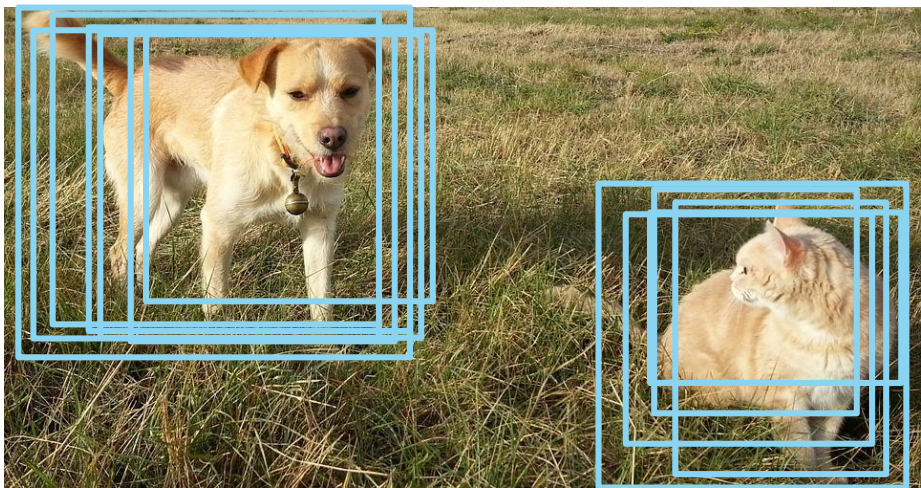"Trying to directly regress to constitutes a difficult learning task"

*Gidaris, Spyros, and Nikos Komodakis. "Locnet: Improving localization accuracy for object detection."*

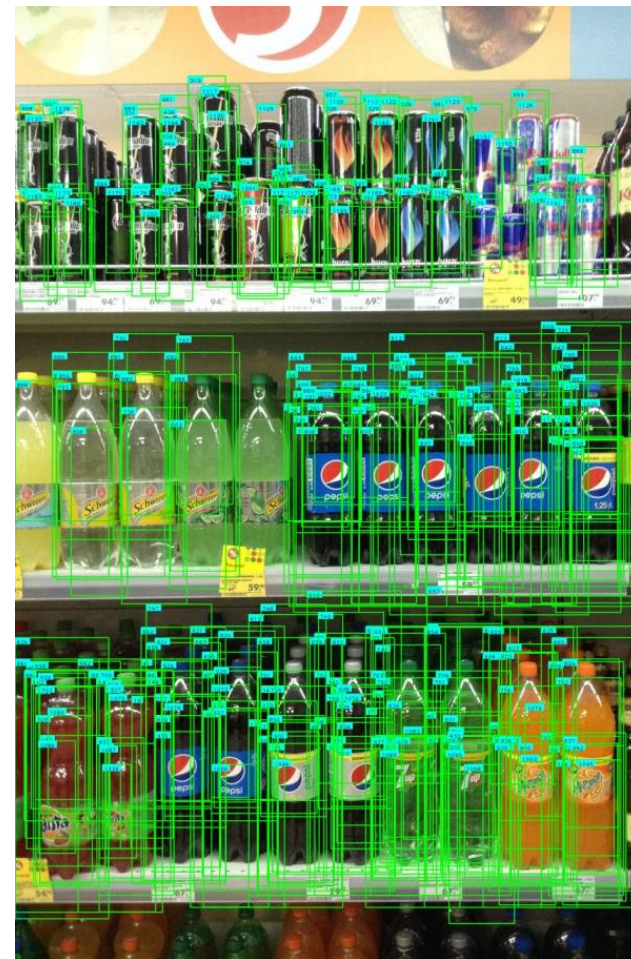"If we pick better priors for the network to start with we can make it easier"

*Redmon, Joseph, and Ali Farhadi. "YOLO9000: Better, Faster, Stronger."*

**Trax** image recognition

# Duplicate Merger

The "Standard" Case

# Are these products  cans?

Yes!

Not Really

No!

Really Not

Trax
image recognition
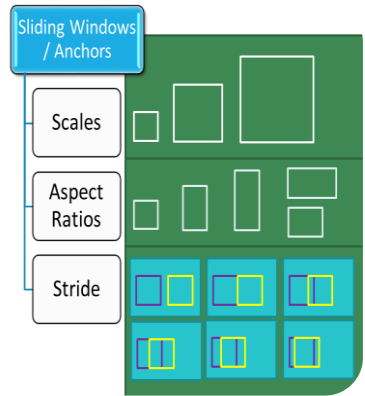
# Detector Innovations - Poster



Region Proposal Network

Objectness Soft Labels

Robust Duplicate Merger

O(1) Region Classification

**Trax** image recognition

# *AGENDA*

❯ **Trax Visual Challenges**

❯ **Deep Context Embedding Architecture**

❯ **Implementation Details**

❯ **The Detection Challenge**

❯ **Summary**

**Trax** image recognition

# *Take Home Message*

> **Fine-grained Classification is Challenging**

> **New Context Embedding CNN Architecture**

> **NLP Inspired**

> **Detection challenge**

**Trax**
image recognition

# Trax
## image recognition

Thank You