

Unreasonable Effectiveness of OCR in Visual Advertisement Understanding

Mayu Otani, Yuki Iwazaki, Kota Yamaguchi
CyberAgent, Inc.



Challenges in visual ads



Symbolism



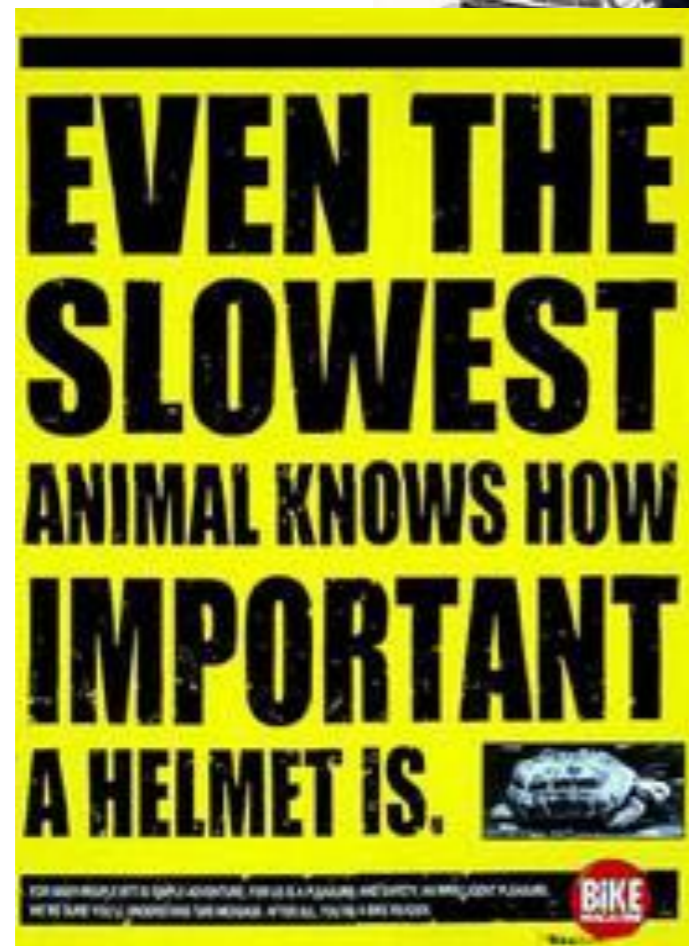
Illustrations



Complicated Compositions

Key observations

Texts on ads seem the most powerful clues



Our approach

1. Extract image2text and ocr2text relevance
2. Rank the statement given both scores

Statement text

Visual and OCR clues from ad

I should shop for jeans at this store
Because they use their profits to build homes for people



Image2text
relevance

old jeans new
homes

Ocr2text
relevance

Our model

"I should shop for jeans ..."



ocr2text network

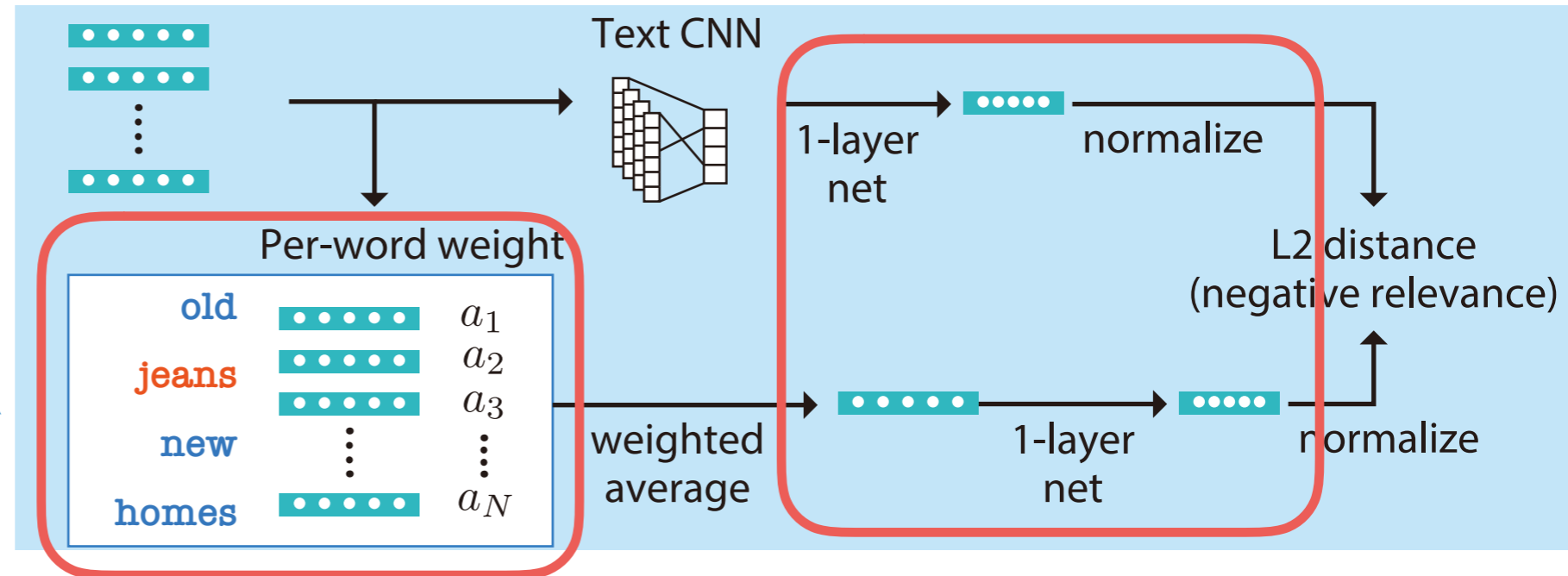
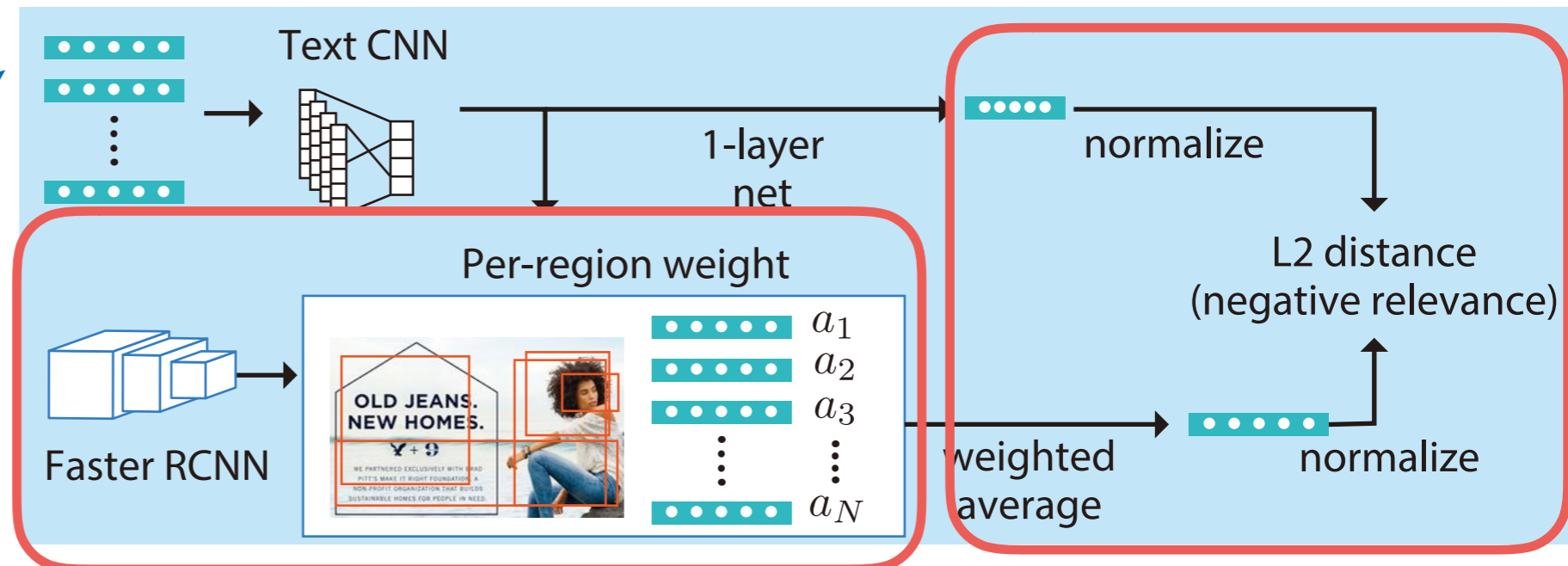
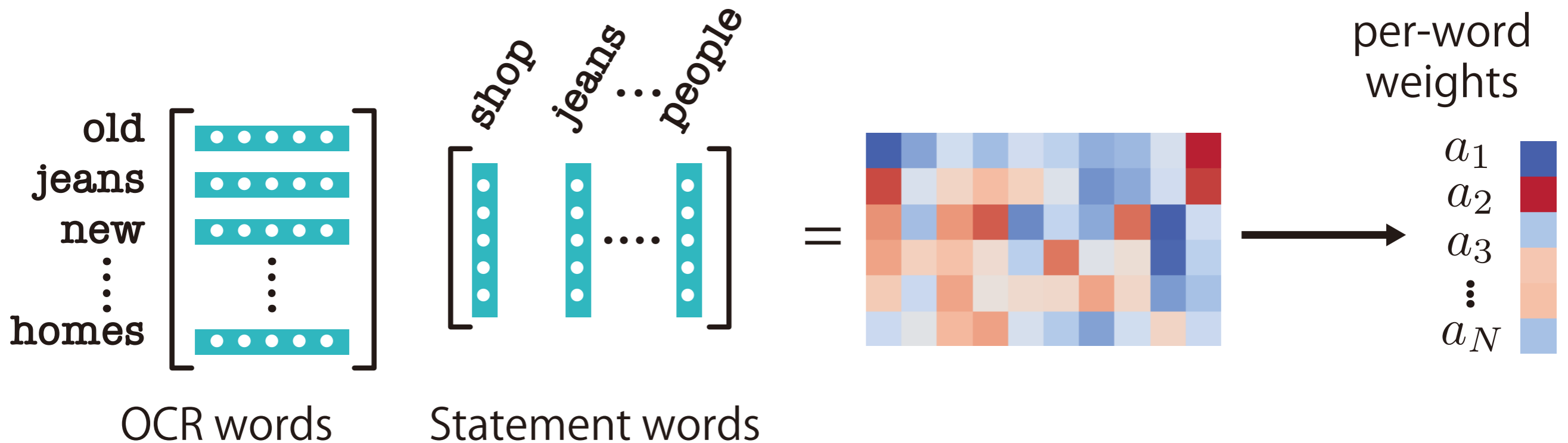


image2text network



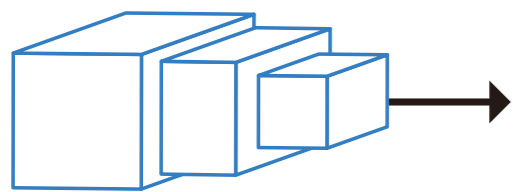
Per-word weights



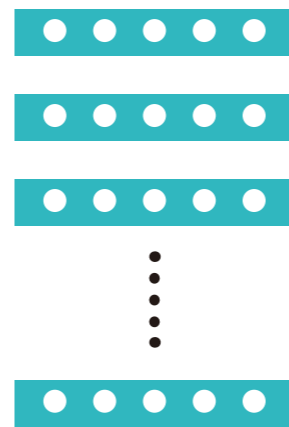
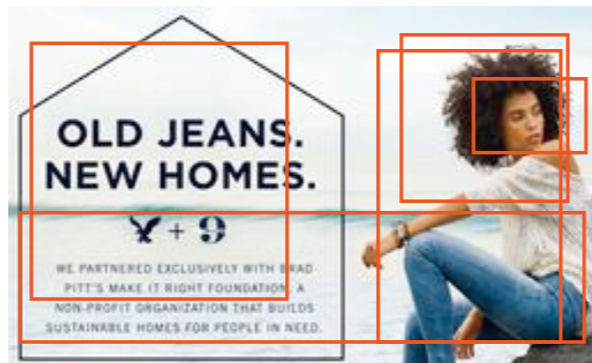
Vector similarity between OCR detected words and words in a statement text

Image embedding module

I should shop for jeans...



Faster RCNN



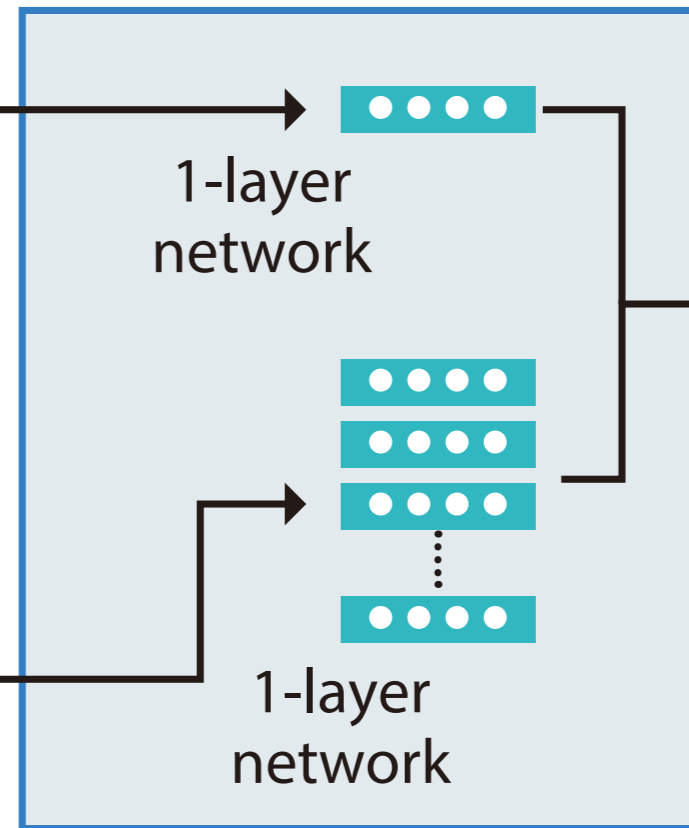
\times

a_{r_1}
 a_{r_2}
 a_{r_3}
 \vdots
 a_{r_M}

Weighted average



Attention network



Per-region weights



Text pre-processing

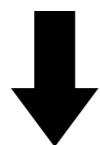
Long-term dependency harms. We trim texts.

I should shop for jeans at this store Because they use their profits to build homes for people

I should <action> because <reason>



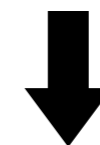
shop for jeans at this store



Action to image / OCR
relevance estimation



they use their profits to build
homes for people



Reason to image / OCR
relevance estimation

Training

Contrastive Loss

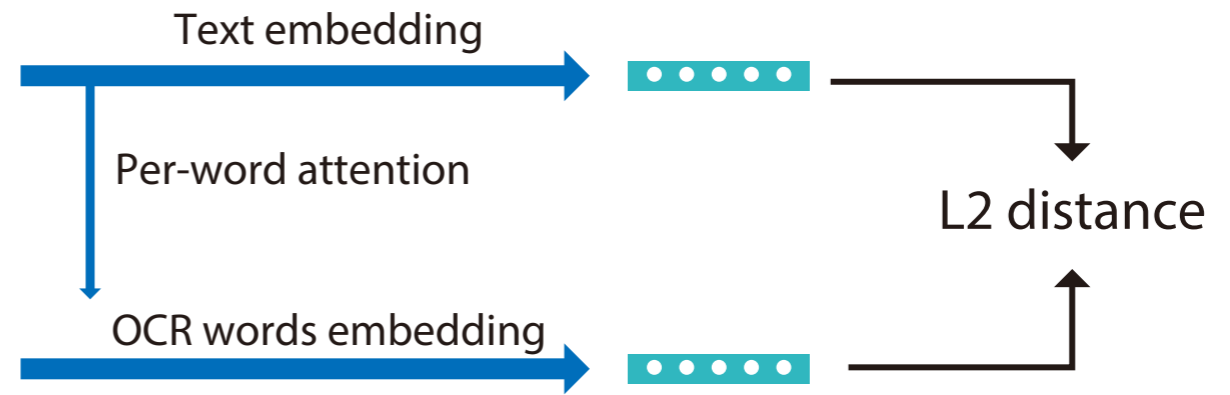
$$L_n = \frac{1}{2} (y_n d_n^2 + (1 - y_n) \max(\mu - d_n, 0)^2)$$

y_n : Label. 1=correct text, 0=incorrect text

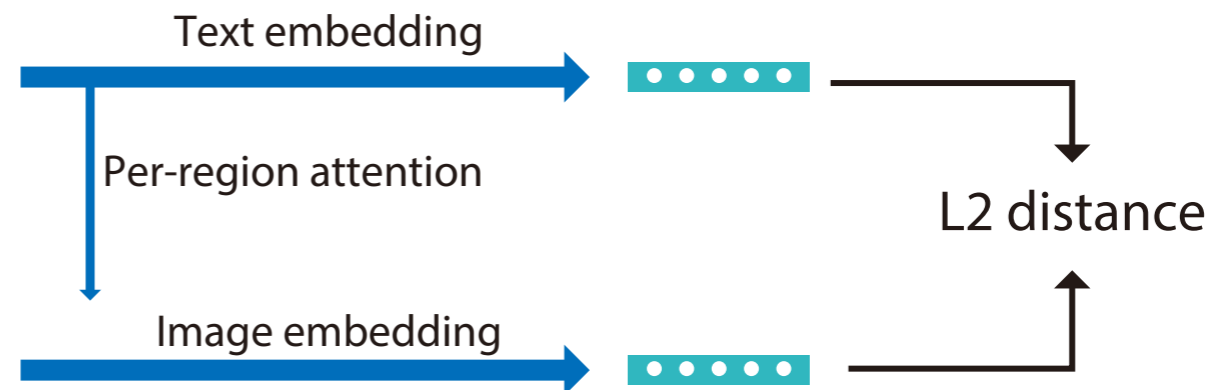
d_n : Distance between text and image/OCR embeddings

“I should shop for jeans ...”

old, jeans, new, homes, ...



“I should shop for jeans ...”



Examples of QA results



A. I should shop for jeans at this store Because they use their profits to build homes for people

Per-word weights for action

old|jeans|new|homes|we|partnered|exclusively|with|brad|make|it|right|foundation|non|profit|organization|that|builds|sustainable|homes|for|people|in|need

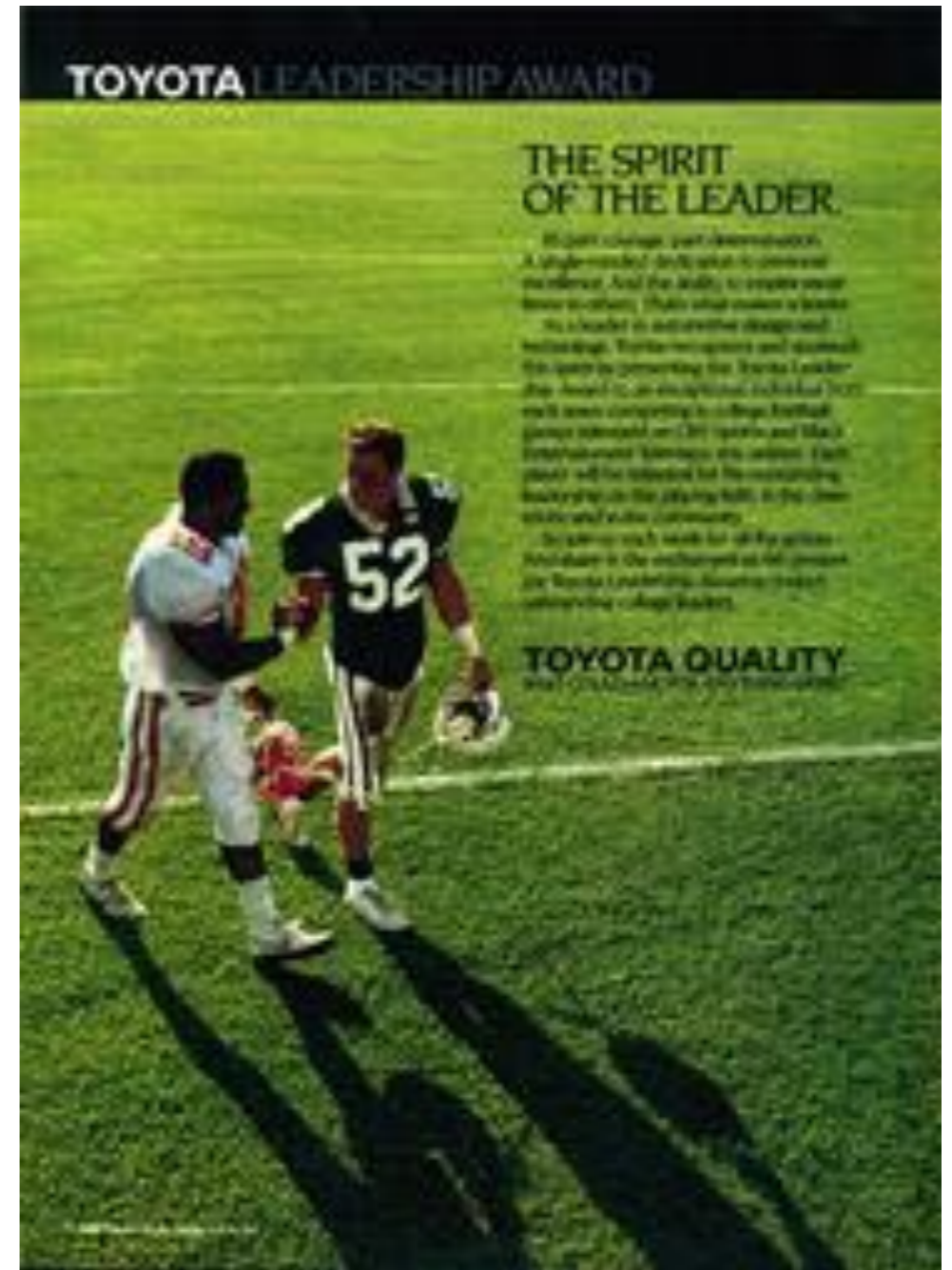
Per-word weights for reason

old|jeans|new|homes|we|partnered|exclusively|with|brad|make|it|right|foundation|non|profit|organization|that|builds|sustainable|homes|for|people|in|need

Examples of QA results

Per-word weights for reason

toyota leadership award the spirit of the
leader part courage part determination
single minded dedication to personal
excellence and the ability to inspire
excel in others what makes leader as
leader in automotive design and
technology toyota recognizes and this
spirit by presenting the toyota leader
ship award to an each team competing
in college football games ...



A. I should buy Toyota Because they support college leaders

Examples of QA results

No words detected



A. I should shop at the gap Because it will make me sexy

How much does OCR help?

P@1

I should shop for jeans at this store Because they use their profits to build homes for people



0.42

I should shop for jeans at this store Because they use their profits to build homes for people



0.82

old jeans
new homes

How much does appearance help?

P@1
(validation set)

I should shop for jeans at this store Because they use their profits to build homes for people



old jeans
new homes

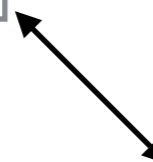


0.83

I should shop for jeans at this store Because they use their profits to build homes for people



0.85

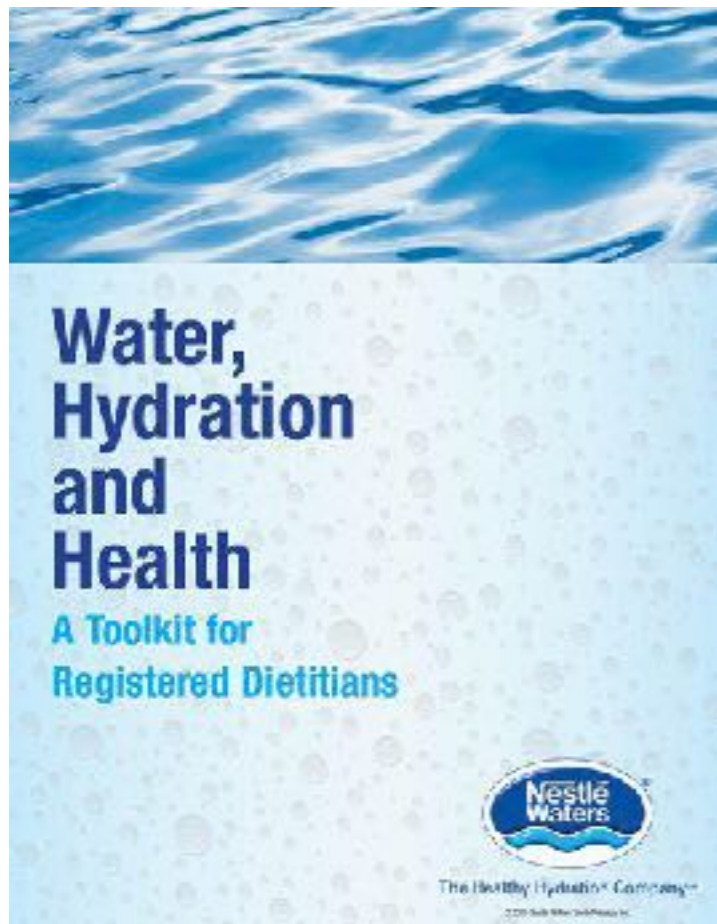


old jeans
new homes

Is OCR quality crucial?

A. Yes

	Score on validation set
Google Cloud Vision OCR	0.85
Tesseract OCR v3	0.54



Tesseract v3

`, 1'2. .4, f4 Water, HydraIiOn and"`



Google Cloud Vision

Water,Hydration and Health
A Toolkit for Registered Dietitians
Nestle Waters
The Healthy Hydration Company

Things we didn't try

but look important

OCR results correction

B L A C K B E R R Y \n R E M E M B E R \n ...



BLACKBERRY REMEMBER ...

Multilingual support

Words detected but discarded



'誕生。美容オイル生恋靴のルージュ。 \nそれは、美容オイルがスティックになった、贅沢な色艶。 \nとろけて密着。色っぽくうるんだ官能の唇へ。 \nマキアージュドラマティックルージュ全10色新発売 \nMAKE uP DATA :pクアチイッタルー問" RD425 /トカルーアイPipp-V1233:般定カラート \n2Mn合わ ■0120.30-47100900-2100庫楽ギ始收定AMHEMO \nwww.shiseido.co.jp mq \nレディにしあがれ。 \nNEWMAQUILLAGE \nとの唇、女っぽくて、ごめんなさい。 \n'

Limitations

- Multilingual support
- Decorative font
- Paintings, drawings



Misc findings

- Higher resolution helps
 - Perhaps because CNNs can literally read texts
 - Also object detection?
- Annotators really read texts
 - And perhaps logos

Symbolism in Ad



Save Water ... Save Life

Save Water ... Save Life

Symbolism might be used for rhetorics, but not necessarily for surface messages.

Remaining questions

- How well do humans understand ad messages without language clues?
- Effects of designs (colors, layout, font style, etc.) in message telling
- Better task / dataset design?
 - Hiding / blurring texts enough?