



### 1. Image Annotation

- Image classification is a challenging task requiring additional information to correctly annotate images.
- We blend visual features extracted from neighbors and their metadata.
- Several convolutional and recurrent neural networks (CNNs-RNNs) are jointly adopted to infer similarity among neighbors and query images.



### 2. Our approach

For a query image  $x_{\!\scriptscriptstyle r}$  we (nonparametrically) generate a neighborhood  $\mathcal{Z}_x$  using metadata and train our neural networks to classify x given  $\mathcal{Z}_x$  . The set of candidate neighborhoods for an image x is the set:

$$\mathcal{Z}_x = \{s \in \mathcal{P}(X): |s| = m\}$$

The prediction s(x, heta) is the average of  $f(x,ec{z}; heta)$  over all candidate neighbourhoods:

$$s(x, heta) = rac{1}{|Z_x|}\sum_{z\in\mathcal{Z}_x}f(x,ec{z}; heta)$$

Our models are trained minimizing the following loss function  ${\cal L}$  :

$$heta^* = arg \min_{ heta} \sum_{(x,y) \in D_{train}} \mathcal{L}(s(x, heta),y)$$

### 3. Metadata Encoding

To correctly recover similar images, our models use metadata which are directly fed to the final layers of NNs after a transformation step.

• One-hot encoding: sum of one-hot vectors for all relevant tags of the query image. Neighborhoods are computed using the Jaccard distance  ${\cal J}$  between binary vectors.

$$p_x = \sum_{i \ s.t. \ t_i \in \{t_{(1)}, t_{(2)}, ..., t_{(n)}\}} e_i^{ au}$$

 Semantic-aware encoding: a transformation step is employed to map binary vectors to a meaningful semantic space. Word2vec and WordNet embeddings are considered.

$$ho(o_x;eta) = \sum_{i=1}^ au o_{x_{(i)}} \, \cdot eta(t_{(i)})$$

# **A CNN-RNN Framework for Image Annotation** from Visual Cues and Social Network Metadata

Tobia Tesan<sup>1</sup>, Pasquale Coscia<sup>2</sup>, Lamberto Ballan<sup>2</sup> <sup>1</sup>Quantexa Ltd, London, UK <sup>2</sup>Department of Mathematics, University of Padova, Italy

## 4. Visual Models

Query images are classified combining only visual features extracted by the AlexNet CNN pre-trained on ImageNet.









## 5. Joint Models

Joint models take advantage of additional information (tags) that is fed to the classification layer after a transformation step.

















### LTwin+2RNN





LTN+AllVecs

LTwin+RNN



Our experiments are performed on the NUS-WIDE dataset considering ~190,000 images and 5000 tags. Per-image/per-label mean Average Precisions (mAPs) metrics show that LTwin model achieves SOTA results compared to several baselines and all the models outperform the visual-only baseline.

Method	mAP <sub>lab</sub>	$mAP_{img}$	reclab	prec <sub>lab</sub>	rec <sub>img</sub>	prec <sub>img</sub>
Tag-only Model + linear SVM [7]	46.67	-	-	-		
Graphical Model (all metadata) [7]	49.00				10 <del>70</del>	
CNN + WARP [16]	-	9 <b>00</b> 0	35.60	31.65	60.49	48.59
CNN-RNN [21]	-	-	30.40	40.50	61.70	49.90
SR-RNN [22]	-		50.17 ×	55.65 <b>*</b>	71.35 <b>*</b>	70.57 <b>*</b>
SR-RNN + Vecs [22] †	-	9 <b>—</b> 3	58.52 <b>*</b>	63.51 <b>*</b>	77.33 *	76.21 <b>*</b>
SRN [35]	60.00	80.60	41.50 <b>*</b>	70.40 <b>*</b>	58.70 <b>*</b>	81.10 *
MangoNet [33]	62.80	80.80	41.00 *	73.90 <b>*</b>	59.90 <b>*</b>	80.60 *
LTN [2]	$52.78 \pm 0.34$	$80.34 \pm 0.07$	$43.61 \pm 0.47$	$46.98 \pm 1.01$	$74.72 \pm 0.16$	$53.69 \pm 0.13$
LTN + Vecs [2] †	$61.88 \pm 0.36$	$80.27 \hspace{0.1in} \pm 0.08$	$57.30 \pm 0.44$	$54.74 \hspace{0.1in} \pm 0.63$	$75.10{\scriptstyle~\pm 0.20}$	$53.46 \pm 0.09$
Upper bound	$100.00 \pm 0.00$	$100.00 \pm 0.00$	$65.82 \pm 0.35$	$60.68 \pm 1.32$	$92.09 \pm 0.10$	$66.83 \hspace{0.1cm} \pm 0.12$
Our baseline: v-only	$45.05 \pm 0.11$	$76.88 \pm 0.11$	$42.31 \pm 0.59$	$43.74 \pm 1.07$	$71.41 \pm 0.13$	$51.36 \pm 0.13$
Our baseline: LTN <sub>n:id</sub>	$53.17 \pm 0.12$	$79.82 \pm 0.16$	$45.67 \pm 1.75$	$47.64 \pm 2.18$	$74.29 \pm 0.13$	$53.34 \pm 0.17$
Our baseline: LTN + Vecs <sub>n:id,f:id</sub> †	$54.86 \pm 0.20$	$81.34 \hspace{0.1cm} \pm 0.15$	$46.56 \pm 1.39$	$50.10 \pm 1.70$	$75.67 \pm 0.17$	$54.37 \pm 0.14$
Our model: RTN <sub>n:w2v</sub>	$55.36 \pm 0.34$	$79.77 \pm 0.27$	$48.73 \pm 2.77$	$51.21 \pm 2.61$	$74.35 \pm 0.29$	$53.28 \pm 0.24$
Our model: LTwin <sub>n:w2v,f:w2v</sub> †	$63.13 \pm 0.31$	$83.77 \pm 0.06$	$54.40 \pm 1.33$	$51.86 \ \pm 1.58$	$78.06 \pm 0.05$	$55.78 \hspace{0.1cm} \pm 0.13$









### 6. Experimental Results

### Visual Models

### Joint Models

E-mail: tobiatesan@quantexa.com, {pasquale.coscia, lamberto.ballan}@unipd.it