



# Scene Text Visual Question Answering

Ali Furkan Biten\*<sup>1</sup>, Rubèn Tito\*<sup>1</sup>, Andrés Mafla\*<sup>1</sup>, Lluis Gomez<sup>1</sup>, Marçal Rusiñol<sup>1</sup>, Ernest Valveny<sup>1</sup>, C.V. Jawahar<sup>2</sup>, Dimosthenis Karatzas<sup>1</sup> <sup>2</sup>CVIT, KCIS, IIIT Hyderabad, India <sup>1</sup>Computer Vision Center, UAB, Spain





Q: What is the price displayed in large letters on the sign? **A:** 14.99



Q: What is written on the sign?



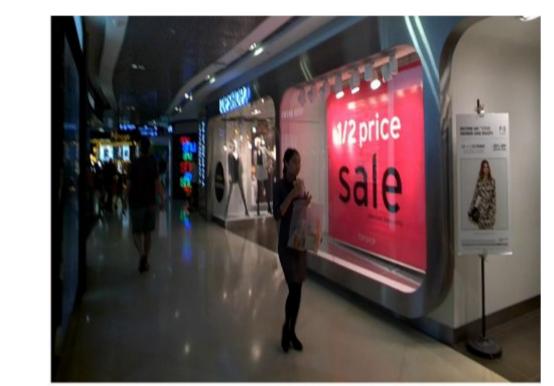
Q: Where is this train going? A: To New York A: New York



Q: What is the exit number on the street sign? A: Exit 2

ST-VQA

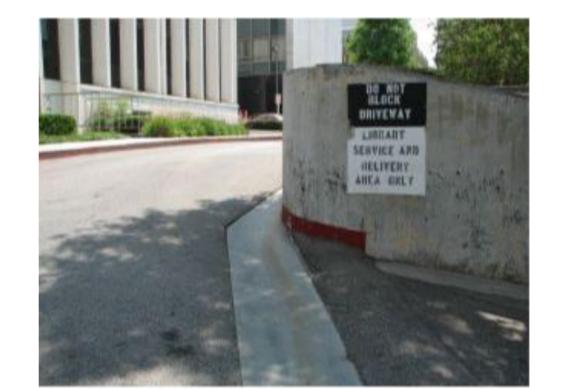
Dataset and Leaderboard



Q: What word in black comes below 1/2 price?



Q: What company's logo is on the A: STARBUCKS COFFEE



Q: What is written in the black A: Do not block driveway

METHOD A

**METHOD E** 

What is the wattage of the vacuum

**Strongly Contextualized Dictionary** 

(100 words per image)

cleaner?

1400W



than the other?



one of the signs is pointing to? served at this establishment? A: Lee Wee Nam Library



Q: Where is the high court located



A: The Ten Commandments



Q: What is the automobile sponsor of the event?

## Overview

- Motivation: Current VQA datasets do not include rich semantic information conveyed as text in an image when reasoning about an image/question pair.
- Images are taken from several commonly used computer vision datasets.
- Question/Answer pairs were collected with Amazon Mechanical Turk.
- New metric is proposed that smoothly captures OCR precision as well as reasoning capabilities.
- A competition and three novel tasks are proposed with increasing levels of complexity.
- Several baselines are proposed and evaluated as a starting point, demonstrating the need of a generative model.

#### ST-VQA Gathering and Statistics

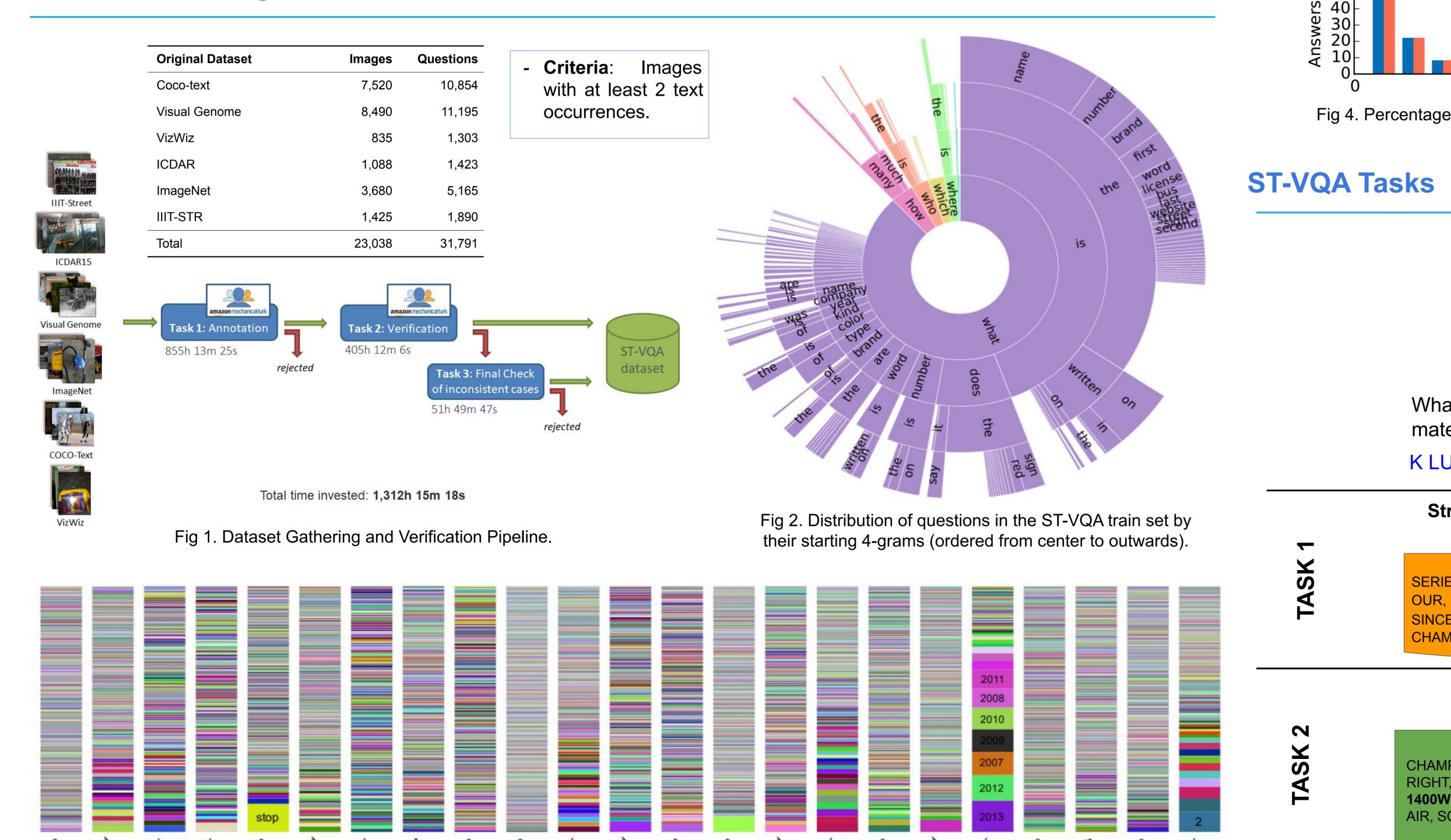
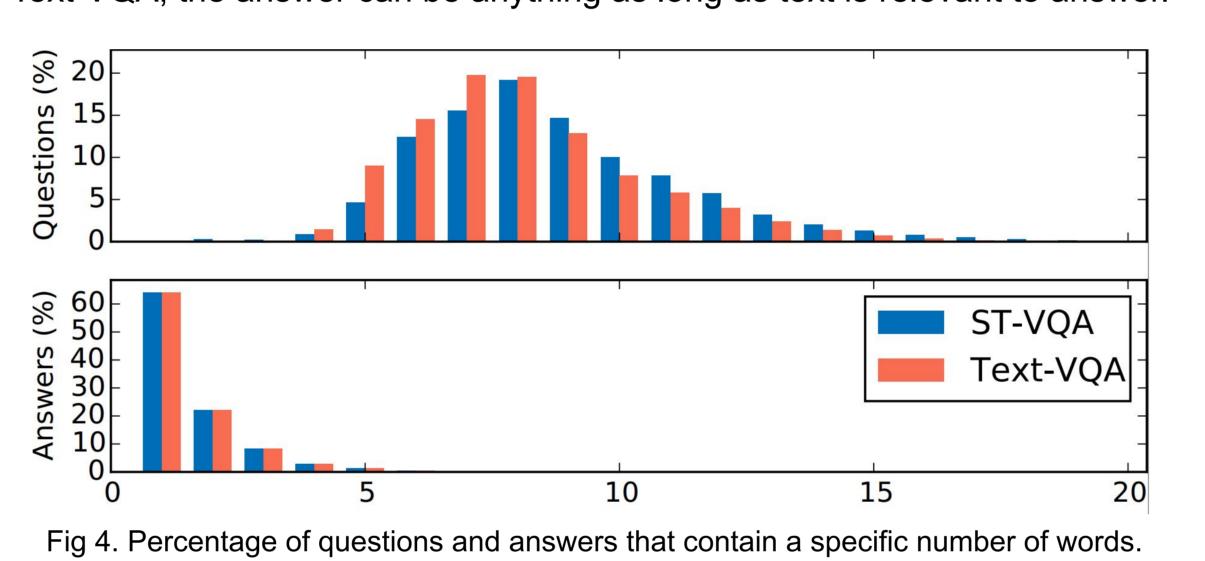


Figure 3. Distribution of answers for different types of questions in the ST-VQA train set. Each color represents a different unique answer.

### **Comparison with TextVQA**

- The two datasets are highly complementary as the image sources used do not intersect with each other.
- We explicitly required a minimum amount of two text instances to be present, while in TextVQA images were sampled on a category basis.
- The answers in ST-VQA are always the text found in the image, whereas in Text-VQA, the answer can be anything as long as text is relevant to answer.



What is written on the piece of

Strongly Contextualized Dictionary

(100 words per image)

material above to walking men?

AMPIONSHIP, MARCH, AROUND,

K LUNC

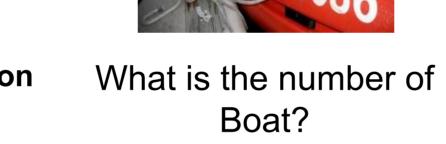
## Metric: Average Normalized Levenshtein Similarity (ANLS)

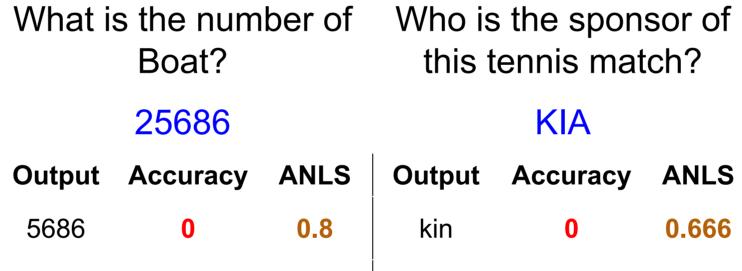
- Tracks accuracy.
- Based on string edit distance.
- Deals with border cases (Correct reasoning but partially wrong recognition).

A: Guinness

$$NLS = \frac{1}{N} \sum_{i=0}^{N} \left( \max_{j} s\left(a_{ij}, o_{qi}\right) \right) \quad \text{, where} \quad s\left(a_{ij}, o_{qi}\right) = \begin{cases} \left(1 - NL\left(a_{ij}, o_{qi}\right)\right) & \text{if } NL\left(a_{ij}, o_{qi}\right) < \tau \\ 0 & \text{if } NL\left(a_{ij}, o_{qi}\right) \geq \tau \end{cases}$$







How do the Omnibuses stop?

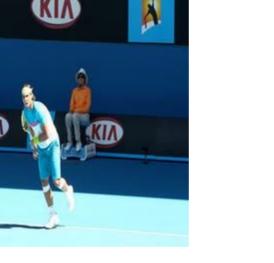
Strongly Contextualized Dictionary

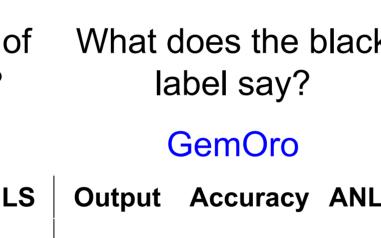
(100 words per image)

JNITED, BUILT, BEGAN, REQUEST, DUE, LONG

VERNMENT, WATER, YORK,

BY REQUEST





**ONE WAY** 

What did the sign say originally?

**Strongly Contextualized Dictionary** 

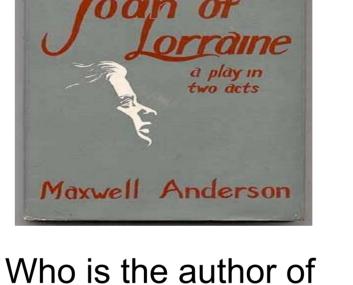
(100 words per image)

ONE, RIGHT, FORCE, OFTEN, FAMILY, POINT,

NATIONAL, BACK, **WAY**, BUILT, AREA,

BEGAN, DUE, LONG, FRANCE, MUCH, PARTY, LIFE

HAMPIONSHIP, MARCH, .



this book? Maxwell Anderson Accuracy ANLS Output Accuracy ANLS Output Accuracy ANLS

## **Q:** What brand are the machines? A: bongard STR (bbox): 1 Scene Image OCR: zbongard **SAAA+STR**: aeropostale SAN(CNN): ray VTA: boar

SAAA+STR: delhi SAN(CNN): delhi VTA: delhi court **USTB-TQA:** delhi QAQ: delhi Clova Al OCR: facal facal

A: delhi

STR (bbox): delhi

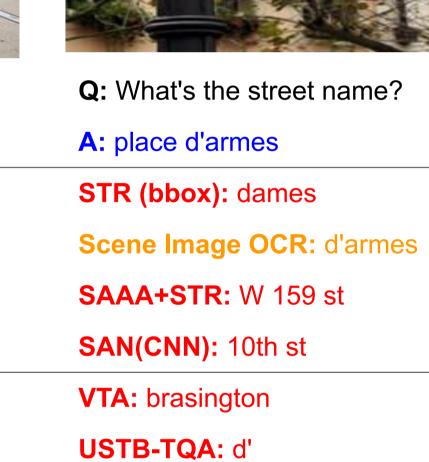
Scene Image OCR: high

**Baseline Qualitative Results** 



A: purple route STR (bbox): purple Scene Image OCR: 1208 SAAA+STR: crosstown **SAN(CNN)**: 66 VTA: purple route **USTB-TQA:** purple QAQ: north

Clova Al OCR: arts- arts-



QAQ: d'armes

Clova Al OCR: place d'armes

# **Baseline Results**

Clova Al OCR: bongard

**USTB-TQA:** boar

**QAQ:** bongard

Method	OCR		V	Task 1		Task 2		Task 3		Upper bound	
		Q		ANLS	Acc.	ANLS	Acc.	ANLS	Acc.	ANLS	Acc.
Random	X	×	X	0.015	0.96	0.001	0.00	0.00	0.00	-	-
STR (retrieval)	<b>✓</b>	X	X	0.171	13.78	0.073	5.55	-	-	0.782	68.84
STR (bbox)	<b>✓</b>	×	X	0.130	7.32	0.118	6.89	0.128	7.21	-	-
Scene Text Image OCR	<b>✓</b>	×	X	0.145	8.89	0.132	8.69	0.140	8.60	-	-
SAAA (1k cls)	X	<b>✓</b>	<b>✓</b>	0.085	6.36	0.085	6.36	0.085	6.36	0.571	31.96
SAAA+STR (1k cls)	<b>✓</b>	<b>✓</b>	<b>✓</b>	0.091	6.66	0.091	6.66	0.091	6.66	0.571	31.96
SAAA (5k cls)	X	<b>✓</b>	<b>✓</b>	0.087	6.66	0.087	6.66	0.087	6.66	0.740	41.03
SAAA+STR (5k cls)	<b>✓</b>	<b>✓</b>	<b>✓</b>	0.096	7.41	0.096	7.41	0.096	7.41	0.740	41.03
SAAA (19k cls)	X	<b>✓</b>	<b>✓</b>	0.084	6.13	0.084	6.13	0.084	6.13	0.862	52.31
SAAA+STR (19k cls)	<b>✓</b>	<b>✓</b>	<b>✓</b>	0.087	6.36	0.087	6.36	0.087	6.36	0.862	52.31
QA+STR (19k cls)	<b>✓</b>	<b>✓</b>	X	0.069	4.65	0.069	4.65	0.069	4.65	0.862	52.31
SAN(LSTM)	X	<b>✓</b>	<b>✓</b>	0.102	7.78	0.102	7.78	0.102	7.78	0.740	41.03
SAN(LSTM)+STR (5k cls)	<b>✓</b>	<b>✓</b>	<b>✓</b>	0.136	10.34	0.136	10.34	0.136	10.34	0.740	41.03
SAN(CNN)+STR (5k cls)	<b>✓</b>	<b>V</b>	<b>✓</b>	0.135	10.46	0.135	10.46	0.135	10.46	0.740	41.03

#### Weakly Contextualized Dictionary (30,000 words, same for all test set images)

ECOND, WILL, SOUTH, FOUR, TRAIN, HIGH, AMERICAN, BRITISH, LONG, OMNIRUSES, FOUR, ANOTHER, LONG, KING, TYPE, WATER, DUE, DUE, GOVERNMENT, ELBA, RIVER, **BY**, PROGRAM, PLACE, BUILDING, AIR, SOUTH, BEGAN, LARGE, NOW, WORK, END, USE, SEVERAL,

> Open Dictionary (no words provided)

MEMBER, 0, PROGRAM,

#### **ICDAR ST-VQA Competition Results**

			Task 1			Task 2					Task 3				
Method	Shared	Specific	Total	Acc.	Dict.	Shared	Specific	Total	Acc.	Dict.	Shared	Specific	Total	Acc.	
VTA	0.507	0.501	0.506	43.52	100	0.280	0.268	0.279	17.77	48.91	0.280	0.285	0.282	18.13	
USTB-TQA	0.457	0.445	0.445	39.98	97.05	0.168	0.196	0.173	13.34	84.11	0.168	0.183	0.170	13.14	
USTB-TVQA	0.129	0.100	0.124	10.09	20.55	0.093	0.094	0.093	6.59	83.76	0.093	0.108	0.095	6.86	
Focus	0.300	0.275	0.295	24.45	68.84	0.080	0.081	0.080	4.16	58.84	0.088	0.089	0.088	4.42	
VQA-DML	0.142	0.138	0.141	11.63	99.97	-	-	-	-	-	-	-	-	-	
TMT	0.076	0.045	0.055	4.53	13.80	-	-	-	-	-	-	-	-	-	
QAQ	-	-	-	-	-	-	-	-	-	-	0.255	0.265	0.256	19.19	
Clova Al OCR	-	-	-	-	-	_	_	-	-	-	0.213	0.224	0.215	12.53	

Biten, A. F., Tito, R., Mafla, A., Gomez, L., Rusiñol, M., Mathew, M., Jawahar, CV, Valveny, E. & Karatzas, D. (2019). ICDAR 2019 Competition on Scene Text Visual Question Answering. 15th IAPR International Conference on Document Analysis and Recognition (ICDAR)