

A SURVEY OF CURRENT DATASETS FOR VISION & LANGUAGE RESEARCH

Microsoft Research

Francis Ferraro¹, Nasrin Mostafazadeh², Ting-Hao (Kenneth) Huang³, Lucy Vanderwende, Jacob Devlin, Michel Galley, Margaret Mitchell

¹ Johns Hopkins University, ² University of Rochester, ³ Carnegie Mellon University

MOTIVATION

Recent explosion in vision&language work, from captioning and video description to question answering and beyond.

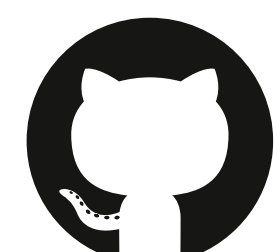
What is lacking? A systematic way to objectively analyze the quality of the datasets:

- Language Quality Criteria
- Vision Quality Criteria

VISIONANDLANGUAGE.NET

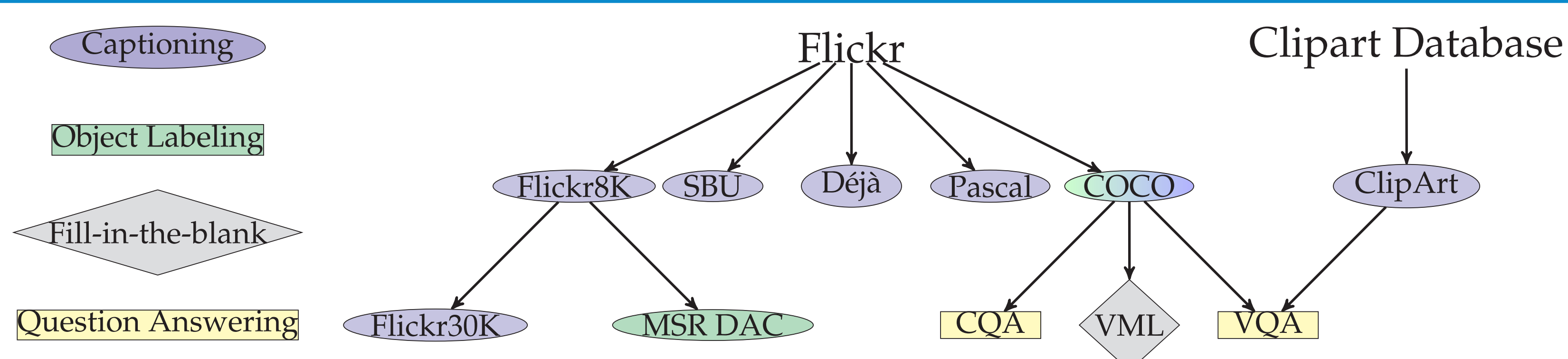
Our evolving website, <http://visionandlanguage.net>, contains pointers and references to a extensive list of older and newly expanding vision & language datasets:

- Categorizes evolving research tasks introduced in vision & language.
- Makes record of the major available datasets and their characteristics.



Source on GitHub, allowing for community involvement in updating.

CATEGORIZATION OF RECENT V&L DATASETS



Flickr8K [1]; Flickr30K [2]; SBU [3]; MSR Dense Annotation Corpus [4]; Déjà Images [5]; Pascal [6]; Microsoft Common Objects in Context [7]; Toronto COCO-Question Answering [8]; Visual MadLibs [9]; Visual Question Answering [10]; Abstract Scenes ClipArt [11]

REPORTING BIAS

Reporting bias and Photographer's bias: People are selective of what they "report" (Gordon & Van Durme, 2013; Torralba & Efros, 2011).

Case Study: Dense Annotation Corpus

(Yatskar et al., 2014)

- 8.04 visible objects per image.
- Only 2.7 (34%) are appear in a caption.

Many selection biases present in abstract scenes are also present in photos.

FINE-GRAINED PERPLEXITY

Test	Brown	Clipart	Coco	CQA	Flickr30K Train	Pascal	SBU	VDC	VQA
VQA	425.9	368.8	366.8	317.7	665.8	119.3	281.0	455.0	19.6
VDC	200.5	52.4	61.5	289.9	51.1	28.7	154.5	30.0	180.1
SBU	473.9	107.1	346.4	328.5	344.0	78.2	119.8	230.7	194.3
Pascal	265.2	64.5	43.2	174.2	63.4	36.0	105.3	83.0	228.2
Flickr30K	247.8	78.5	54.3	181.5	37.8	39.9	125.0	72.1	192.2
CQA	489.4	186.1	137.0	33.8	244.5	74.9	200.1	259.0	72.1
Coco	274.6	59.2	36.2	137.0	75.3	39.3	111.0	87.1	236.9
Clipart	233.6	11.2	117.4	210.8	109.4	28.7	130.6	82.5	114.7
Brown	237.1	99.6	560.8	354.0	405.0	47.8	621.5	187.3	126.5

Some datasets (COCO, Flickr30K, Clipart) may be more useful as out-of-domain data.

LANGUAGE QUALITY

The distinction between a qualitatively "good" or "bad" dataset is task dependent; all criteria should be viewed through the lens of particular downstream tasks.

- **Vocabulary Size**
- **Average Sentence Length**
- **Part of Speech Distribution**
- **Syntactic Complexity (Yngve, Frazier):** measures embedding/branching in a sentence's syntax.
- **Abstract/Concrete Ratio:** Percentage of abstract terms (from Vanderwende et al. (2015)) indicates the range of visual and non-visual concepts the dataset covers.
- **Perplexity:** Compared against 5-gram language model learned on generic 30B words English.

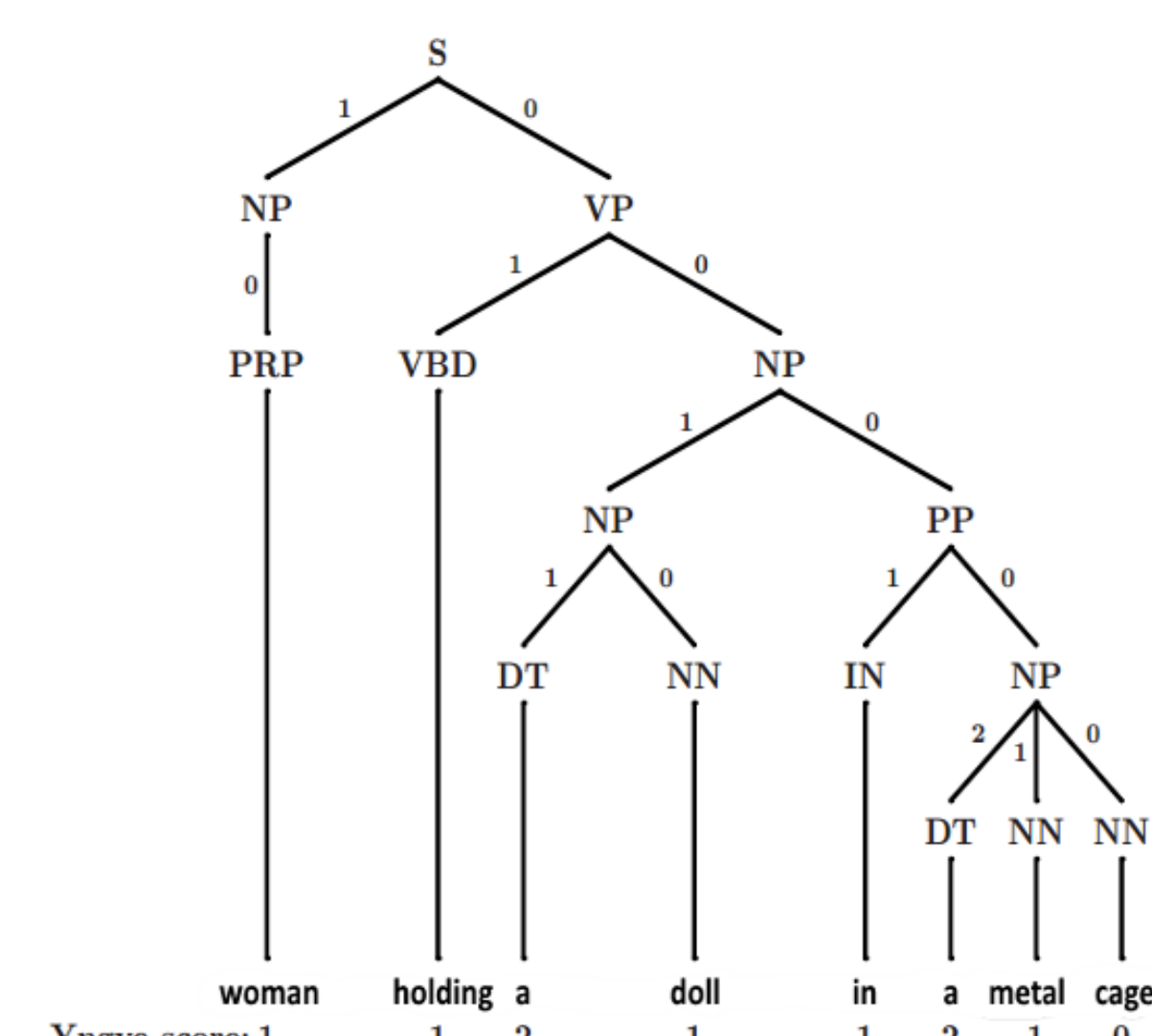
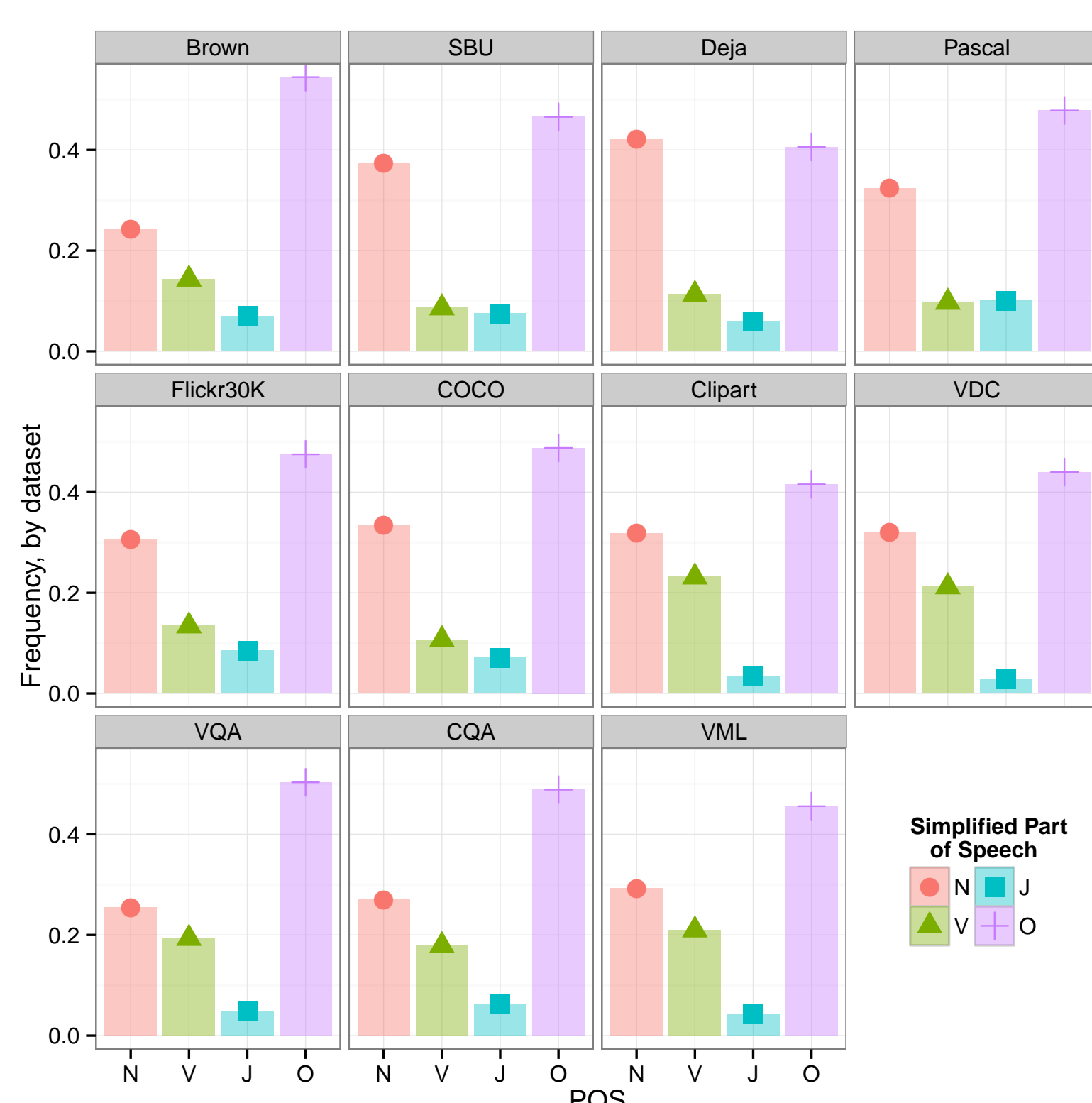


Figure: Example parse tree with branch scores for Yngve scoring.

POS DISTRIBUTION



	Dataset	Size(k)		Language					
		Img	Txt	Frazier	Yngve	Vocab Size (k)	Sent Len.	%Abs	Ppl
Balanced	Brown	-	52	18.5	77.21	47.7	20.82	15.24%	194
	SBU	1000	1000	9.70	26.03	254.6	13.29	3.74%	346
User-Gen	Déjà	4000	180	4.13	4.71	38.3	4.10	9.70%	184
	Pascal	1	5	8.03	25.78	3.4	10.78	17.74%	123
Crowd-sourced	Flickr30K	32	159	9.50	27.00	20.3	12.98	14.98%	118
	COCO	328	2500	9.11	24.92	24.9	11.30	12.96%	121
	Clipart	10	60	6.50	12.24	2.7	7.18	17.96%	126
Video	VDC	2	85	6.71	15.18	13.6	7.97	12.86%	148
	VQA	10	330	6.50	14.00	6.2	7.58	19.22%	113
Beyond	CQA	123	118	9.69	11.18	10.2	8.65	16.14%	199
	VML	11	360	6.83	12.72	11.2	7.56	17.19%	110

Vision Quality Criteria in the Paper!

REFERENCES

1. Rashtchian et al., (2010)
2. Young et al., (2014)
3. Ordonez et al., (2011)
4. Yatskar et al., (2014)
5. Chen et al., (2015)
6. Farhadi et al., (2010)
7. Lin et al., (2014)
8. Ren et al., (2015)
9. Yu et al., (2015)
10. Antol et al., (2015)
11. Zitnick et al., (2013)