

Paraphrase Acquisition from Image Captions

EACL 2023



**Marcel
Gohsen¹**



**Matthias
Hagen²**



**Martin
Potthast³**

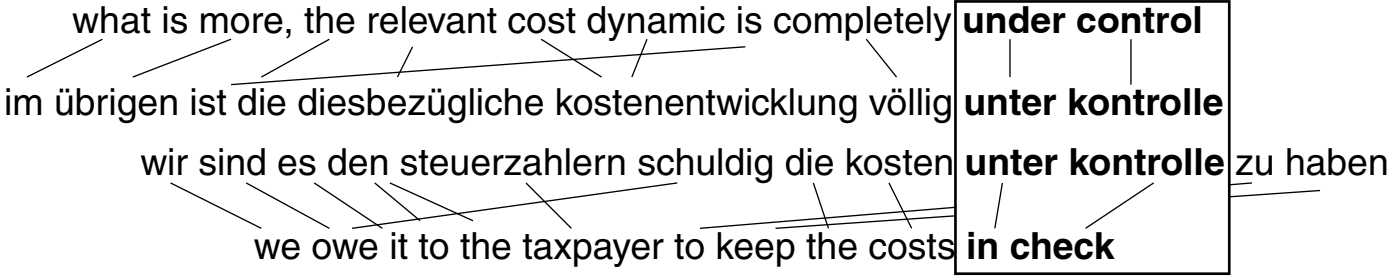


**Benno
Stein¹**



Pivoting for Paraphrase Acquisition

Pivoting is a common distantly-supervised method to acquire paraphrases.



“Paraphrasing with Bilingual Parallel Corpora”, Bannard and Callison-Burch 2005

Account
@usesname

Samsung halts oriduction of its **Galaxy Note 7** as battery problems linger.

9:30 PM • Frb 5, 2022

Account
@usesname

Samsung temporarily supended production of its **Galaxy Note 7** devices following reports.

9:30 PM • Frb 5, 2022

“A Continuosly Growing Dataset of Sentential Paraphrases”, Lan et al. 2017

Image Captions as Paraphrases



WIKIPEDIA
The Free Encyclopedia

Easter Bunny

From Wikipedia, the free encyclopedia



Ēostre

From Wikipedia, the free encyclopedia



Image Captions as Paraphrases



WIKIPEDIA
The Free Encyclopedia

Easter Bunny

From Wikipedia, the free encyclopedia



A 1907 postcard featuring
the Easter Bunny

Ēostre

From Wikipedia, the free encyclopedia



An Easter postcard from
1907 depicting a rabbit

Image Captions as Paraphrases

Infobox captions

Easter Bunny

🌐 48 languages ▾

Article Talk

Read View source View history Tools ▾

From Wikipedia, the free encyclopedia



The **Easter Bunny** (also called the **Easter Rabbit** or **Easter Hare**) is a folkloric figure and symbol of **Easter**, depicted as a **rabbit**—sometimes dressed with clothes—bringing **Easter eggs**. Originating among German **Lutherans**, the "Easter Hare" originally played the role of a judge, evaluating whether children were good or disobedient in behavior at the start of the season of **Eastertide**,^[1] similar to the "naughty or nice" list made by **Santa Claus**. As part of the legend, the creature carries colored eggs in its basket, as well as candy, and sometimes toys, to the homes of children. As such, the Easter Bunny again shows similarities to Santa (or the **Christkind**) and **Christmas** by bringing gifts to children on the night before a holiday. The custom was first^t^[2]^[unreliable source?] mentioned in **Georg Franck von Franckenuau's** *De ovis paschalibus*^[3] ('About Easter eggs') in 1682, referring to a German tradition of an Easter Hare bringing eggs for the children.

Symbols

Rabbits and hares

The hare was a popular motif in medieval church art. In ancient times, it was widely believed (as by **Pliny**, **Plutarch**, **Philostratus**, and **Aelian**) that the hare was a **hermaphrodite**.^[4]^[5]^[6] The idea that a hare could reproduce without loss of **virginity** led to an association with the **Virgin Mary**, with hares sometimes occurring in **illuminated manuscripts** and **Northern European** paintings of the **Virgin** and **Christ Child**. It may also have been associated with the **Holy Trinity**, as in the **three hares** motif.^[4]^[7]^[unreliable source?]^[8]

Eggs

Main articles: *Easter egg* and *Egg decorating*

Eggs have been used as **fertility** symbols since **antiquity**.^[9] Eggs became a symbol in Christianity associated with rebirth, especially the 1st century AD, with the resurrection of the

Easter Bunny



A 1907 postcard featuring the Easter Bunny

Grouping	Legendary creature
Sub grouping	Animal
Other name(s)	Easter Rabbit, Easter Hare
Country	Germany



Image and Caption Mining

Image Filter

Caption Filter

Image Clustering

Paraphrase Construction

Paraphrase Filter



Wikipedia-IPC

Image Captions as Paraphrases

Imagebox captions

Connection to Easter Hares

In Northern Europe, Easter imagery often involves [hares](#) and [rabbits](#).^[31] The first scholar to make a connection between the goddess Eostre and hares was Adolf Holtzmann in his book *Deutsche Mythologie*. Holtzmann wrote of the tradition, "the Easter Hare is inexplicable to me, but probably the hare was the sacred animal of Ostara; just as there is a hare on the statue of [Abnoba](#)." Citing folk [Easter customs](#) in [Leicestershire](#), England, where "the profits of the land called Harecrop Leys were applied to providing a meal which was thrown on the ground at the 'Hare-pie Bank'", late 19th-century scholar [Charles Isaac Elton](#) speculated on a connection between these customs and the worship of [Ēostre](#).^[32] In his late 19th-century study of the hare in folk custom and mythology, Charles J. Billson cited numerous incidents of folk customs involving hares around the Easter season in Northern Europe. Billson said that "whether there was a goddess named [Ēostre](#), or not, and whatever connection the hare may have had with the ritual of Saxon or British worship, there are good grounds for believing that the sacredness of this animal reaches back into an age still more remote, where it is probably a very important part of the great Spring Festival of the prehistoric inhabitants of this island."^[22]

Adolf Holtzmann had also speculated that "the hare must once have been a bird, because it lays eggs" in modern German folklore. From this statement, numerous later sources built a modern legend in which the goddess Eostre transformed a bird into an egg-laying hare.^[33] A response to a question about the origins of Easter hares in the 8 June 1889 issue of the journal *American Notes and Queries* stated: "In Germany and among the Pennsylvania Germans toy rabbits or hares made of canton flannel stuffed with cotton are given as gifts on Easter morning. The children are told that this Osh'ter has laid the Easter eggs. This curious idea is thus explained: The hare was originally a bird, and was changed into a quadruped by the goddess Ostara; in gratitude to Ostara or Eastre, the hare exercises its original bird function to lay eggs for the goddess on her festal day."^[34] According to folklorist Stephen Winick, by 1900, many popular sources had picked up the story of Eostre and the hare. One described the story as one of the oldest in mythology, "despite the fact that it was then less than twenty years old."^[33]

Some scholars have further linked customs and imagery involving hares to both [Ēostre](#) and the Norse goddess [Freyja](#). Writing in 1972, John Andrew Boyle cited commentary contained within an etymology dictionary by A. Ernout and A. Meillet, where the authors write that "Little else ... is known about [[Ēostre](#)], but it has been suggested that her lights, as goddess of the dawn, were carried by hares. And she certainly represented spring [fecundity](#), and love and carnal pleasure that leads to fecundity." Boyle responded that nothing is known about [Ēostre](#) outside of Bede's single passage, that the authors had seemingly accepted the identification of [Ēostre](#) with the Norse

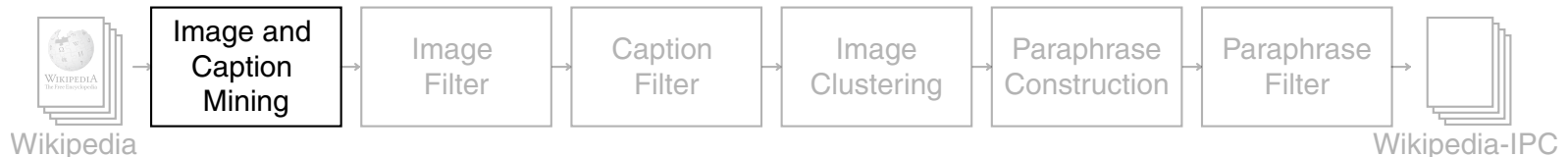


Image Captions as Paraphrases

Alternative texts



Riihimäki, Finland 2009

Slovak Easter symbols

Easter Bunny postcard, 1907

Easter Bunny postcard, 1907



★ *Oryctolagus cuniculus*

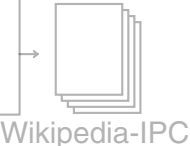
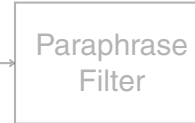
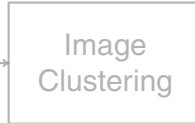
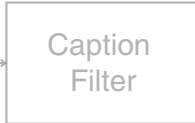
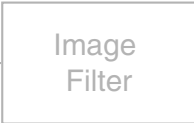


Image Captions as Paraphrases

Goal: Discard icons, symbols and pictograms.

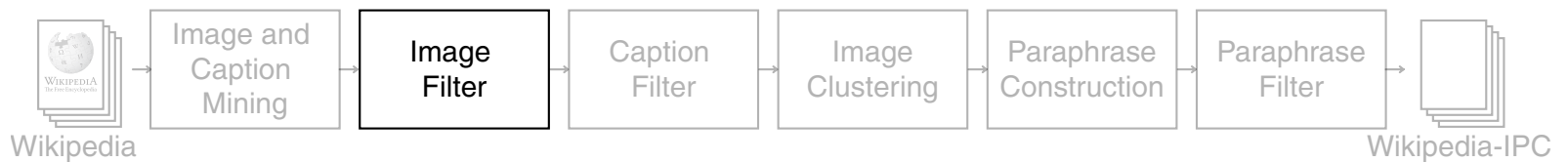


Image Captions as Paraphrases

Goal: Discard icons, symbols and pictograms.



Icons, symbols and pictograms appear much more frequently in Wikipedia articles.

Men's sponsored sports by school [edit]

School	Baseball	Basketball	Cross Country	Football	Golf	Tennis	Track & Field (Indoor)	Track & Field (Outdoor)	Total Southland Sports
Houston Christian	✓	✓	✓	✓	✓	✗	✓	✓	7
Incarnate Word	✓	✓	✓	✓	✓	✓	✓	✓	8
Lamar	✓	✓	✓	✓	✓	✓	✓	✓	8
McNeese	✓	✓	✓	✓	✓	✗	✓	✓	7
New Orleans	✓	✓	✓	✗	✓	✓	✓	✓	7
Nicholls	✓	✓	✓	✓	✓	✓	✗	✗	6
Northwestern State	✓	✓	✓	✓	✗	✗	✓	✓	6
Southeastern Louisiana	✓	✓	✓	✓	✓	✗	✓	✓	7
Texas A&M-Commerce	✗	✓	✓	✓	✓	✗	✓	✓	6
Texas A&M-Corpus Christi	✓	✓	✓	✗	✗	✓	✓	✓	6
Totals	9	10	10	8	8	5	9	9	67

⇒ Limit number of references to [2, 10]

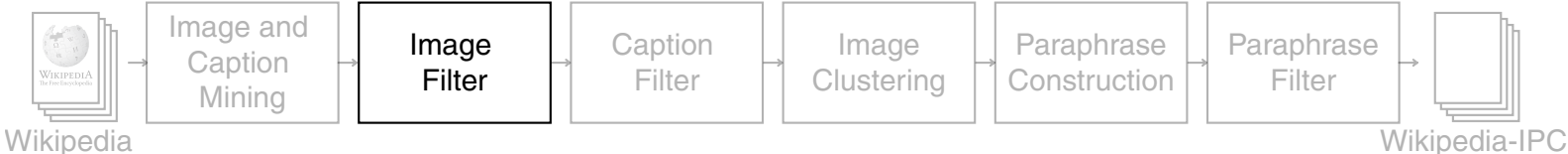


Image Captions as Paraphrases

Goal: Discard “trivial” image captions.

- Too short captions (<6 words)

“Easter Bunny Postcard, 1907”

“Altar in temple”

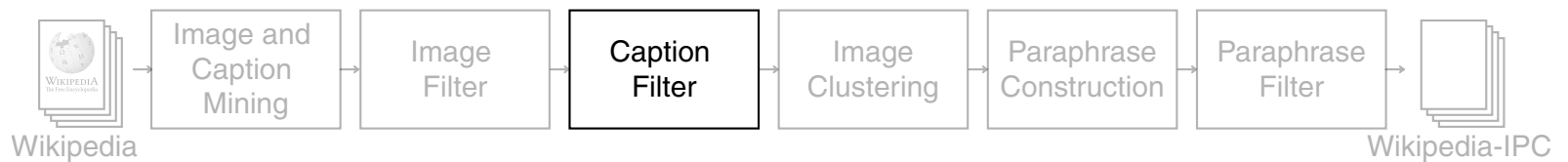


Image Captions as Paraphrases

Goal: Discard “trivial” image captions.

- ❑ Too short captions (<6 words)

“Easter Bunny Postcard, 1907”

“Altar in temple”

- ❑ Non-sentential captions (noun phrases)

“Ted Danson, Best Actor in a Comedy Series winner”

“Players in Grant Park during Pokemon Go Fest 2017”

“Postcard of the large rectangular Grand Central Station”

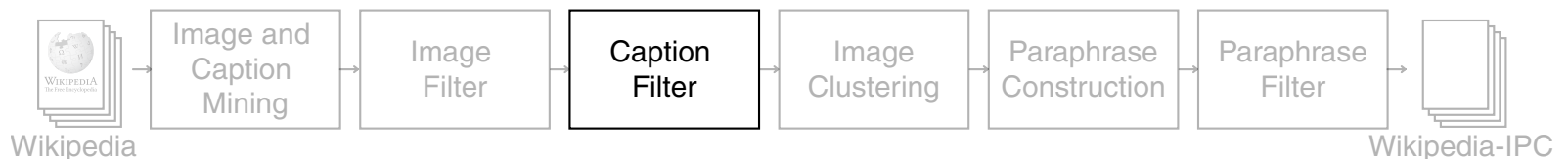
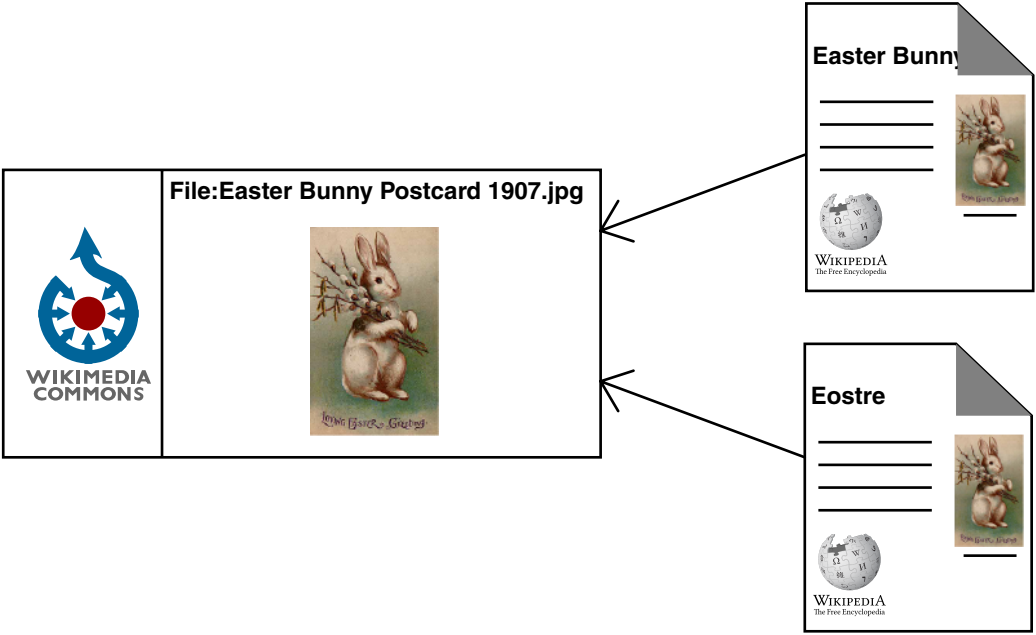


Image Captions as Paraphrases

Goal: Cluster equivalent images.



- Usage on en.wikipedia.org
 - [Easter Bunny](#)
 - [Eostre](#)
 - [User talk:Airplaneman/Archive 10](#)
 - [User:Stebunik](#)
 - [User:Dr. Blofeld/April 2015](#)

⇒ Image url is a sufficient equivalence criterion

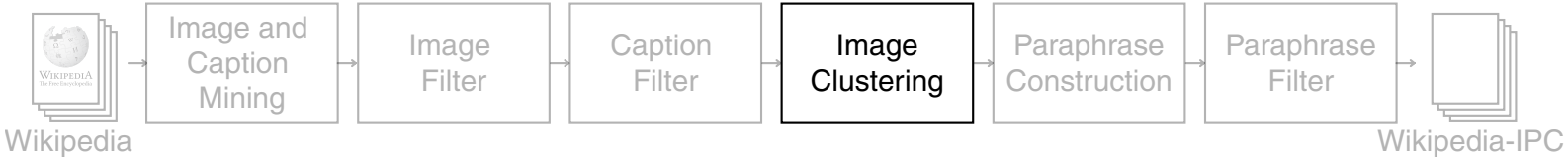



Image Captions as Paraphrases

Image Cluster



Captions

- A 1907 postcard featuring the Easter Bunny
- An Easter postcard from 1907 depicting a rabbit


Alternative texts

- A hare standing on its hind legs and carrying several branches
- A drawing of an Easter bunny carrying several branches as part of an Easter postcard



Image Captions as Paraphrases

Image Cluster



Captions

A 1907 postcard featuring the Easter Bunny

An Easter postcard from 1907 depicting a rabbit

A 1907 postcard featuring the Easter Bunny

An Easter postcard from 1907 depicting a rabbit

Alternative texts

A hare standing on its hind legs and carrying several branches

A drawing of an Easter bunny carrying several branches as part of an Easter postcard

A hare standing on its hind legs and carrying several branches

A drawing of an Easter bunny carrying several branches as part of an Easter postcard



Image Captions as Paraphrases

Goal: Discard exact and near duplicates.

Near duplicates are caption pairs that only differ in

- Casing

“A postcard featuring the Easter Bunny” ↔ *“A postcard featuring the **e**aster **b**unny”*

- Punctuation

“A postcard featuring the Easter Bunny” ↔ *“A postcard**,** featuring the Easter Bunny.”*

- Bracket expressions

“A postcard featuring the Easter Bunny” ↔ *“A postcard featuring the Easter Bunny **(1907)**”*

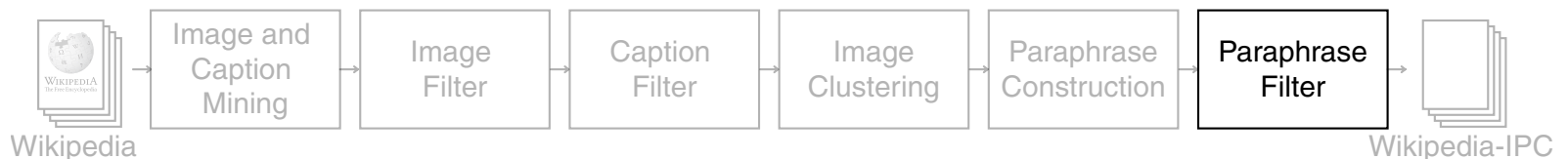


Image Captions as Paraphrases

Gold-quality paraphrases

- Captions pairs from Wikipedia without revision history which are sentences

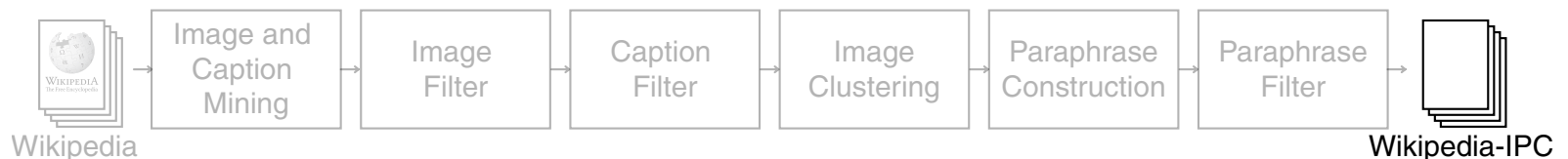
Silver-quality paraphrases

- Gold-quality paraphrases and captions pairs that are sentence fragments

Bronze-quality paraphrases

- Caption pairs from Wikipedia with revision history with sentence fragments

Quality	Caption pairs
Gold	30,237
Silver	229,877
Bronze	656,560



Quantitative Similarity Analysis

Goal: Compare paraphrases from datasets on multiple similarity dimensions.

Lexical and syntactic similarity metrics:

- ❑ ROUGE-1 [Lin, 2004]
- ❑ ROUGE-L [Lin, 2004]
- ❑ BLEU [Papineni et al., 2002]

Quantitative Similarity Analysis

Goal: Compare paraphrases from datasets on multiple similarity dimensions.

Lexical and syntactic similarity metrics:

- ❑ ROUGE-1 [Lin, 2004]
- ❑ ROUGE-L [Lin, 2004]
- ❑ BLEU [Papineni et al., 2002]

Semantic similarity metrics:

- ❑ Word Mover Distance (WMS) [Kusner et al., 2015]
- ❑ BERTScore [Zhang et al., 2019]
- ❑ Sentence Transformer (ST) [Reuners and Gurevych, 2019]

Quantitative Similarity Analysis

$\Delta_{\text{sem,syn}}$

Sophisticated paraphrases are

- semantically as similar as possible
- syntactically as different as possible

$\Delta_{\text{sem,syn}}$ is the average semantic similarity minus average syntactic similarity.

Quantitative Similarity Analysis

$\Delta_{\text{sem},\text{syn}}$

Sophisticated paraphrases are

- semantically as similar as possible
- syntactically as different as possible

$\Delta_{\text{sem},\text{syn}}$ is the average semantic similarity minus average syntactic similarity.

Image Caption Pairs	Syntactic similarity				Semantic similarity				$\Delta_{\text{sem},\text{syn}}$
	ROUGE-1	ROUGE-L	BLEU	Avg.	WMS	BERT	ST	Avg.	
An Easter postcard from 1907 depicting a rabbit. A 1907 postcard featuring the Easter Bunny.	0.53	0.13	0.14	0.27	0.76	0.76	0.90	0.81	0.54

Troops clearing rubble after the May air raid on Belfast. Soldiers clearing rubble after the May air raid on Belfast.	0.90	0.90	0.99	0.93	0.92	0.89	0.98	0.93	0.00

Quantitative Similarity Analysis

Corpus	Acquisition	Syntactic similarity				Semantic similarity				$\Delta_{\text{sem,syn}}$
		ROUGE-1	ROUGE-L	BLEU	Avg.	WMS	BERT	ST	Avg.	
Wikipedia-IPC _{gold}	Caption	0.74	0.71	0.56	0.67	0.83	0.69	0.91	0.81	0.14
Wikipedia-IPC _{silver}	Caption	0.71	0.67	0.52	0.63	0.63	0.81	0.90	0.78	0.15
Flickr8k	Caption	0.53	0.48	0.22	0.41	0.59	0.73	0.86	0.73	0.32
MS-COCO	Caption	0.51	0.45	0.22	0.39	0.57	0.71	0.86	0.71	0.32
PASCAL	Caption	0.51	0.47	0.22	0.40	0.59	0.72	0.86	0.73	0.32
ParaNMT-5m	Generated	0.63	0.60	0.33	0.52	0.60	0.75	0.87	0.74	0.22
PAWS	Generated	0.94	0.79	0.69	0.81	0.82	0.96	0.97	0.92	0.11
MSRPC	Distant supervision	0.73	0.69	0.54	0.65	0.72	0.82	0.90	0.82	0.16
PPDB 2.0	Distant supervision	0.64	0.63	0.32	0.53	0.64	0.63	0.89	0.72	0.19
TaPaCo	Distant supervision	0.65	0.63	0.30	0.53	0.78	0.79	0.91	0.83	0.30

Quantitative Similarity Analysis

Corpus	Acquisition	Syntactic similarity				Semantic similarity				$\Delta_{\text{sem,syn}}$
		ROUGE-1	ROUGE-L	BLEU	Avg.	WMS	BERT	ST	Avg.	
Wikipedia-IPC _{gold}	Caption	0.74	0.71	0.56	0.67	0.83	0.69	0.91	0.81	0.14
Wikipedia-IPC _{silver}	Caption	0.71	0.67	0.52	0.63	0.63	0.81	0.90	0.78	0.15
Flickr8k	Caption	0.53	0.48	0.22	0.41	0.59	0.73	0.86	0.73	0.32
MS-COCO	Caption	0.51	0.45	0.22	0.39	0.57	0.71	0.86	0.71	0.32
PASCAL	Caption	0.51	0.47	0.22	0.40	0.59	0.72	0.86	0.73	0.32
ParaNMT-5m	Generated	0.63	0.60	0.33	0.52	0.60	0.75	0.87	0.74	0.22
PAWS	Generated	0.94	0.79	0.69	0.81	0.82	0.96	0.97	0.92	0.11
MSRPC	Distant supervision	0.73	0.69	0.54	0.65	0.72	0.82	0.90	0.82	0.16
PPDB 2.0	Distant supervision	0.64	0.63	0.32	0.53	0.64	0.63	0.89	0.72	0.19
TaPaCo	Distant supervision	0.65	0.63	0.30	0.53	0.78	0.79	0.91	0.83	0.30

Quantitative Similarity Analysis

Corpus	Acquisition	Syntactic similarity				Semantic similarity				$\Delta_{\text{sem,syn}}$
		ROUGE-1	ROUGE-L	BLEU	Avg.	WMS	BERT	ST	Avg.	
Wikipedia-IPC _{gold}	Caption	0.74	0.71	0.56	0.67	0.83	0.69	0.91	0.81	0.14
Wikipedia-IPC _{silver}	Caption	0.71	0.67	0.52	0.63	0.63	0.81	0.90	0.78	0.15
Flickr8k	Caption	0.53	0.48	0.22	0.41	0.59	0.73	0.86	0.73	0.32
MS-COCO	Caption	0.51	0.45	0.22	0.39	0.57	0.71	0.86	0.71	0.32
PASCAL	Caption	0.51	0.47	0.22	0.40	0.59	0.72	0.86	0.73	0.32
ParaNMT-5m	Generated	0.63	0.60	0.33	0.52	0.60	0.75	0.87	0.74	0.22
PAWS	Generated	0.94	0.79	0.69	0.81	0.82	0.96	0.97	0.92	0.11
MSRPC	Distant supervision	0.73	0.69	0.54	0.65	0.72	0.82	0.90	0.82	0.16
PPDB 2.0	Distant supervision	0.64	0.63	0.32	0.53	0.64	0.63	0.89	0.72	0.19
TaPaCo	Distant supervision	0.65	0.63	0.30	0.53	0.78	0.79	0.91	0.83	0.30

Qualitative Similarity Analysis

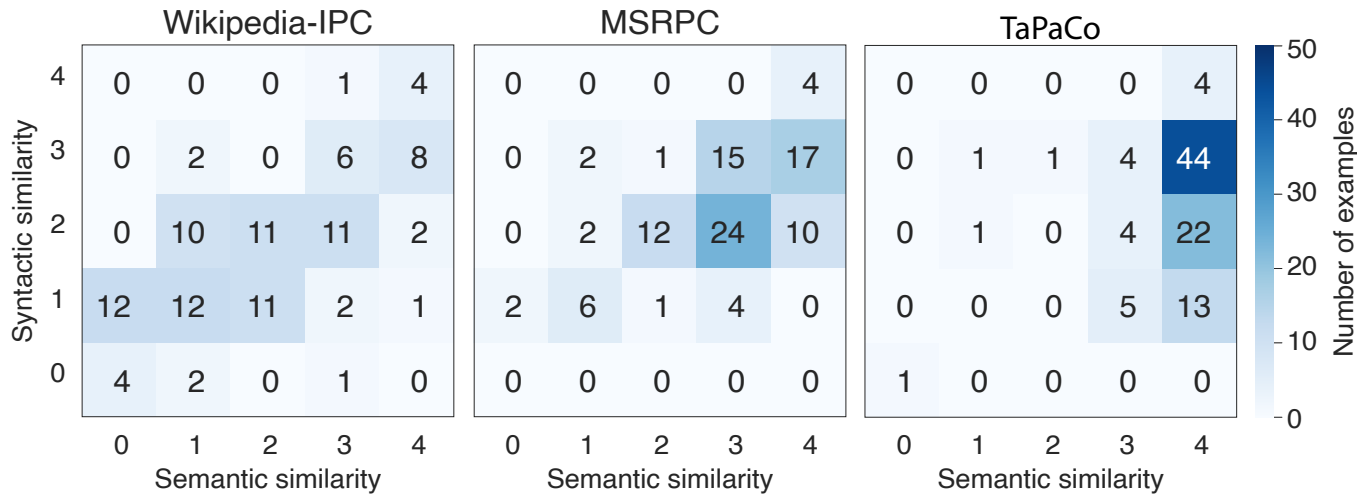
Annotated a paraphrase's semantic and syntactic similarity on a 5-point Likert scale

- ❑ Three datasets (Wikipedia-IPC, MSRPC, TaPaCo)
- ❑ Sample of 100 paraphrase pairs per dataset
- ❑ Four expert annotators

Qualitative Similarity Analysis

Annotated a paraphrase's semantic and syntactic similarity on a 5-point Likert scale

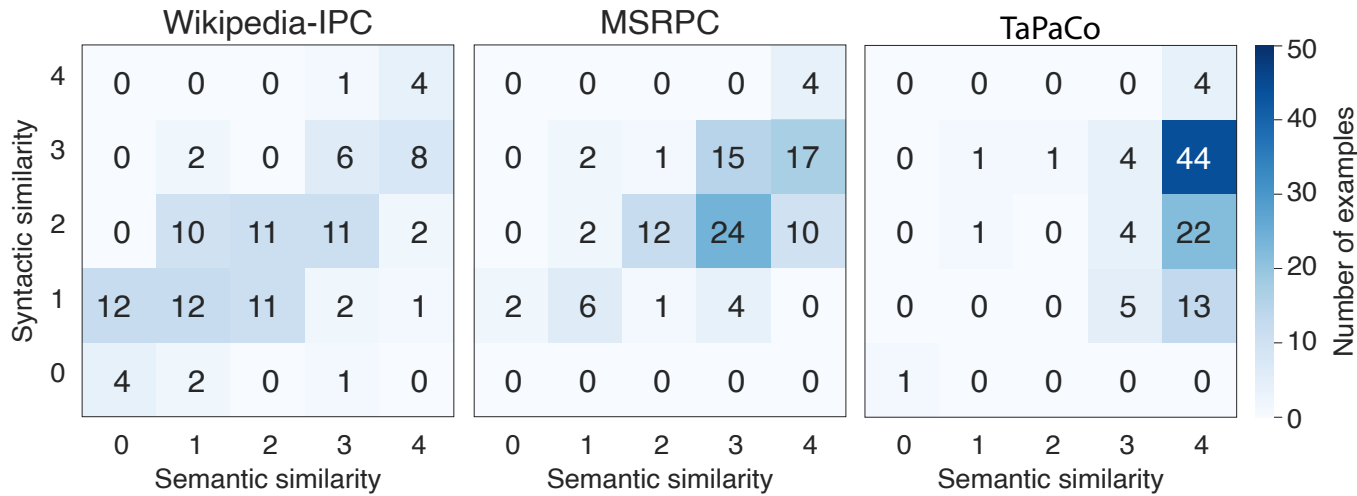
- Three datasets (Wikipedia-IPC, MSRPC, TaPaCo)
- Sample of 100 paraphrase pairs per dataset
- Four expert annotators



Qualitative Similarity Analysis

Annotated a paraphrase's semantic and syntactic similarity on a 5-point Likert scale

- Three datasets (Wikipedia-IPC, MSRPC, TaPaCo)
- Sample of 100 paraphrase pairs per dataset
- Four expert annotators



Syntax:	ROUGE-1	ROUGE-L	BLEU	Average
<i>r</i>	0.78	0.77	0.70	0.79
Semantics:	WMS	BERT	ST	Average
<i>r</i>	0.59	0.70	0.78	0.76

Conclusion

Contributions

- Paraphrase acquisition from image captions
- Paraphrase dataset with different quality levels: Wikipedia-IPC
- Paraphrase sophistication metric: $\Delta_{\text{sem},\text{syn}}$

Conclusion

Contributions

- Paraphrase acquisition from image captions
- Paraphrase dataset with different quality levels: Wikipedia-IPC
- Paraphrase sophistication metric: $\Delta_{\text{sem},\text{syn}}$

Future Work

- Apply paraphrase acquisition method to larger Web resources
- Incorporate image analysis for image equivalence classification

Conclusion

Contributions

- Paraphrase acquisition from image captions
- Paraphrase dataset with different quality levels: Wikipedia-IPC
- Paraphrase sophistication metric: $\Delta_{\text{sem},\text{syn}}$

Future Work

- Apply paraphrase acquisition method to larger Web resources
- Incorporate image analysis for image equivalence classification

Code and Data



<https://github.com/webis-de/EACL-23>

Conclusion

Contributions

- Paraphrase acquisition from image captions
- Paraphrase dataset with different quality levels: Wikipedia-IPC
- Paraphrase sophistication metric: $\Delta_{\text{sem},\text{syn}}$

Future Work

- Apply paraphrase acquisition method to larger Web resources
- Incorporate image analysis for image equivalence classification

Code and Data



Thank you!