

# On the Role of Summary Content Units in Text Summarization Evaluation

Anonymous EMNLP submission

## Abstract

At the heart of the pyramid evaluation method for text summarization lie human written summary content units (SCUs). These SCUs are concise sentences that decompose a summary into small facts. Such SCUs can be used to judge the quality of a candidate summary, possibly partially automated via natural language inference (NLI) systems. Interestingly, with the aim to fully automate the pyramid evaluation, [Zhang and Bansal \(2021\)](#) show that SCUs can be approximated from parsed semantic role triplets (STUs). However, several questions currently lack answers, in particular i) Are there other ways of approximating SCUs that can offer advantages? ii) Under which conditions are SCUs (or their approximations) offering the most value? In this work, we examine two novel strategies to approximate SCUs: generating SCU approximations from AMR meaning representations (SMUs) and from large language generation models (SGUs), respectively. We find that while STUs and SMUs are competitive, the best approximation quality is achieved by SGUs. We also show through a simple sentence-decomposition baseline (SSUs) that SCUs (and their approximations) offer the most value when ranking short summaries, but may not help as much when ranking systems or longer summaries.

## 1 Introduction

Judging the quality of a summary is a challenging task. Besides being short and faithful to its source document, a summary should particularly excel in *relevance*, that is, a summary should select only the most relevant or salient facts from a source document. An attractive method for assessing such notion of relevance is the *Pyramid*-method ([Nenkova and Passonneau, 2004](#)) that is based on so-called *Summary Content Units* (SCUs) which decompose a reference summary into concise human-written English sentences. With SCUs available from one or different reference summaries, we can then more

objectively assess the degree to which a candidate summary contains the relevant information. With the aim to fully automate the Pyramid method, [Zhang and Bansal \(2021\)](#) suggest that the required human effort can be partially and even fully alleviated, by i) automatically generating SCUs and ii) validating the relevance of a summary with an NLI system that checks how many SCUs are entailed by a candidate summary.

Since strong NLI systems are available off-the-shelf and are known to be useful in NLG evaluation<sup>1</sup>, clearly the generation of SCUs is the most challenging and least-understood part of an automated pyramid. Indeed, while [Zhang and Bansal \(2021\)](#) show that SCUs can be approximated by phrasing semantic role parsed triplets, we lack availability and understanding of possible alternatives as well as their potential impact on downstream-task summary evaluation in different scenarios.

In this work, we proposed two novel approaches to approximate SCUs: SMUs that are based on abstract meaning representation (AMR) and SGUs that leverage SoTA large language models (LLMs). We carry out experiments to systematically evaluate the intrinsic quality of SCUs and their approximations. On the downstream task evaluation, we find that although SCUs remain the most effective metric to rank different systems or summaries across three meta-evaluation datasets, surprisingly, an efficient sentence splitting baseline already yields competitive results when compared to SCUs. In fact, the sentence splitting baseline outperforms the best SCU approximation method on a few datasets when ranking systems or long summaries.

In summary, our work provides important insights into the application of automation of the pyramid method in different scenarios for evaluating summaries. We make the code publicly available at [URL Upon Acceptance].

<sup>1</sup>E.g., see [Chen and Eger \(2022\)](#), or [Steen et al. \(2023\)](#).

## 2 Related work

Evaluating the quality of a summary is a challenging task. Over the past two decades, researchers have proposed a wide range of human-in-the-loop or automatic metrics to assess summaries in different dimensions, including linguistic quality, coherence, faithfulness, and content quality. For more in-depth surveys on this topic, please refer to Howcroft et al. (2020) and Gehrmann et al. (2022).

In this work, we focus on evaluating the content quality of a summary that assesses whether the summary effectively captures the salient information of interest from the input document(s). In the reference-based metrics, content quality is assessed by comparing system-generated summaries to human-written reference summaries. The pyramid method (Nenkova and Passonneau, 2004) is regarded as a reliable and objective approach to evaluating the content quality of a summary. Below we briefly describe the pyramid method and highlight some previous efforts to automate this method.

**Pyramid Method.** The original pyramid method (Nenkova and Passonneau, 2004) comprises two steps: SCUs creation and system evaluation. In the first step, human annotators exhaustively identify Summary Content Units (SCUs) from the reference summaries. Each SCU is a concise sentence and describes a single fact. The weight of an SCU is determined by the number of references in which it occurs. In the second step, the presence of each SCU in a system summary is manually checked. The system summary’s pyramid score is calculated as the normalized sum of the weights of the SCUs that are present. Later, Shapira et al. (2019) introduce a revised version of the original pyramid method where they eliminate the merging and weighting of SCUs, thereby enabling SCUs with the same meaning to coexist.

**Automation of the Pyramid Method.** Given the high cost and the expertise required for implementing the pyramid method, in recent years there are a few attempts to automate this approach. Peyrard et al. (2017) propose an automatically learned metric to directly predict human pyramid scores based on a set of features. Zhang and Bansal (2021) propose a system called *Lite<sup>3</sup> Pyramid* that uses a semantic role labeller to extract semantic triplet units (STUs) to approximate SCUs. They further use a trained natural language inference (NLI) model to

replace the manual work of assessing SCUs’ presence in system summaries. In our work, we explore two new methods to approximate SCUs. We further investigate the effectiveness of the automated pyramid method in different scenarios.

## 3 SCU approximation I: SMU from AMR

*Abstract Meaning Representation* (AMR) (Banarescu et al., 2013) is a widely-used semantic formalism employed to encode the meaning of natural language text in the form of rooted, directed, edge-labeled, and leaf-labeled graphs. The AMR graph structure facilitates machine-readable explicit representations of textual meaning.

Motivated by Zhang and Bansal (2021)’s observation that STUs based on semantic roles cannot well present single facts in long reference summary sentences that contain a lot of modifiers, adverbial phrases, or complements, we hypothesize that AMR has the potential to capture such factual information more effectively. This is because, in addition to capturing semantic roles, AMR models finer nuances of semantics, including negations, inverse semantic relations, and coreference.

To generate semantic meaning units (SMUs) from a reference summary, we first leverage a pre-trained Text2AMR model<sup>2</sup> to represent each sentence in the summary as an AMR graph. We then design a few heuristics to split each AMR graph into several sub-graphs and apply an AMR2Text model<sup>3</sup> on each sub-graph to generate SMUs. Please refer to Appendix A.1 for more details on splitting an AMR graph into multiple sub-graphs.

## 4 SCU approximation II: SGU from LLM

Recently, it became widely known that pre-trained large language models (LLMs) are able to generate high-quality output according to prompts given by humans, optionally exploiting shown examples through in-context learning (Brown et al., 2020). Therefore, we try to approximate SCUs using GPT models from OpenAI, calling the resulting units as Semantic GPT Units (SGUs). Specifically, we use GPT-3.5-Turbo which is built on InstructGPT (Ouyang et al., 2022) and GPT-4 (OpenAI, 2023) to generate SGUs (SGUs<sub>3.5</sub> and SGUs<sub>4</sub>) for each reference summary using the same prompt and a

<sup>2</sup>*parse\_xfm\_bart\_large* (<https://github.com/bjascob/amrplib-models>)

<sup>3</sup>*generate\_t5wtense* (<https://github.com/bjascob/amrplib-models>)

one-shot example. Please refer to Appendix A.2 for more details.

## 5 Experiments

### 5.1 Dataset and NLI models

**Data.** We run the experiments on four existing English meta-evaluation datasets: (1) TAC08 (NIST, 2008), (2) TAC09 (NIST, 2009), (3) REALSumm (Bhandari et al., 2020) and (4) PyrXSum (Zhang and Bansal, 2021) and evaluated the results on the last three datasets, using TAC08 for development purposes. Each dataset contains one or multiple reference summaries, the corresponding human-written SCUs, the generated summaries from different systems, and the human evaluation result for each summary/system based on the pyramid method. Table 1 shows some statistics of the reference summaries across different datasets. In general, PyrXSum contains short and abstractive summaries, while RealSumm and TAC09 contain long and extractive summaries. More details on the datasets can be found in appendix A.4.

**NLI Models.** We use the NLI model from Zhang and Bansal (2021) that was fine-tuned on TAC08’s SCU presence gold annotations based on a NLI model from Nie et al. (2020).

### 5.2 Baselines

**STUs** are the semantic role triples based on semantic role labelling (Zhang and Bansal, 2021).

**Sentence splitting** is a baseline that may shed light on the overall usefulness of SCUs in summary evaluation. We split every reference summary into sentences and treat them as SCU approximations.

**N-grams** consist of phrases randomly extracted from a reference summary. For each sentence from the summary, we naïvely generate all possible combinations of 3, 4, and 5 consecutive words. We then randomly select a subset from these combinations, which accounts for 5% of the total number of n-grams produced.

### 5.3 Intrinsic Evaluation

As proposed by Zhang and Bansal (2021), we evaluate approximation quality with an *easiness score*. The score is built by iterating over each SCU-SxU pair and average over the maximum ROUGE-1-F1 score found for each SCU. Naturally this score is recall-biased, and therefore, we also present the

	RealSumm	PyrXSum	TAC09
Avg. # sent.	4.73	2.02	27.22
Avg. # words	63.71	20.56	386.82
Avg. # words/sent	13.47	10.18	14.21
# ref summary	1	1	4
Avg. # SCUs	10.56	4.78	31.63

Table 1: Statistics of the reference summaries from different datasets.

Metrics	RealSumm		PyrXSum		TAC09	
	R	P	R	P	R	P
sentence split	.54	.67	.41	.54	.50	.54
ngrams	.41	.52	.38	.52	.46	.39
STUs	<b>.66</b>	.68	.54	.65	<b>.61</b>	.53
SMUs	.56	.58	.53	.58	.52	.48
SGUs_3.5	.58	.67	.58	.63	.36	.48
SGUs_4	.61	<b>.69</b>	<b>.61</b>	<b>.66</b>	.52	<b>.61</b>

Table 2: Intrinsic Evaluation Results. R is the recall oriented simulation easiness score by (Zhang and Bansal, 2021), while P is our precision-oriented score that is computed in the reverse direction.

score calculated in the reverse direction, to assess whether our approximated SCUs are of high precision.

The results are shown in Table 2. We find that best approximation quality for RealSumm is achieved by STUs, while for PyrXSum, SGU\_4 performs best. Considering the longer texts of TAC09, STUs excel in recall, while SGU\_4 excels in precision.

### 5.4 Extrinsic Evaluation

Our downstream evaluation consists of two parts: summary quality evaluation at the system and summary levels, respectively. System-level correlation assessment evaluates the ability of the metric to compare different summary systems individually. In contrast, summary-level evaluation determines the metric’s ability to compare summaries created by different systems for a common set of documents. Following (Zhang and Bansal, 2021), we use Pearson  $r$  and Spearman  $\rho$  to evaluate the correlations between metrics with gold human labelling scores. Pearson measures linear correlation and Spearman measures ranking correlation. See more details on how to compute these two types of correlation in appendix A.3.

The results are shown in Table 3. We can observe that SGUs overall offer the most useful SCU approximation, with the exception for TAC09

Metrics	System-Level						Summary-Level					
	RealSumm		PyrXSum		TAC09		RealSumm		PyrXSum		TAC09	
	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$
SCUs	.95	.95	.98	.98	.99	.97	.59	.58	.70	.69	.76	.70
SCU Approximations												
- sentence split	.93	<b>.95</b>	<b>.97</b>	<b>.97</b>	.97	.94	.48	.46	.37	.36	<b>.73</b>	.66
- ngrams	.90	.92	.94	.82	.96	.92	.36	.35	.38	.38	.65	.61
- STUs	.92	.94	.95	.95	<b>.98</b>	.95	.51	.50	.46	.44	<b>.73</b>	<b>.67</b>
- SMUs	<b>.94</b>	.94	.96	.94	<b>.98</b>	<b>.96</b>	.50	.48	.46	.44	.70	.64
- SGUs_3.5	.93	<b>.95</b>	<b>.97</b>	.93	.96	.88	.49	.46	.56	.55	.54	.49
- SGUs_4	.92	.94	<b>.97</b>	.95	<b>.98</b>	<b>.96</b>	<b>.54</b>	<b>.52</b>	<b>.58</b>	<b>.56</b>	.71	.66

Table 3: Results of different metrics on three datasets. Best numbers among all SCU approximations are bolded.

(summary-level), where STUs remain the best approximation method, slightly outperforming our simple sentence splitting baseline. However, SGUs still lack the usefulness of true SCUs, which seem to remain the most useful way to evaluate summary quality (if resources permit). Interestingly, however, to discriminate the quality of systems, it is enough to use any approximation, even the sentence split baseline is sufficient to accurately discriminate systems.

## 5.5 Human Evaluation

For a representative sample of human results of our experiment, three authors evaluated the quality of SCUs, STUs, SMUs and SGUs\_4 for 10 reference summaries randomly sampled from REALSumm and PyrXSum. Cohen’s  $\kappa$  scores among three annotators range from 0.37 to 0.87. More details about human evaluation can be found in Appendix A.5.

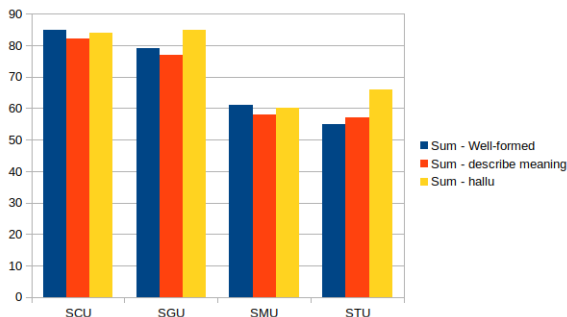


Figure 1: Human evaluation results. Numbers on the y-axis represent the aggregated scores of all three annotators for the 10 examples, with higher scores indicating better performance across all three dimensions.

The findings presented in Figure 1 illustrate that the quality of SMUs is comparable to the STUs. SMUs are a little bit more well-formed but h-

lucinate more compared to STUs. Furthermore, it’s important to note that SGUs are nearly on par with SCUs in terms of overall quality. However, despite their close performance, SGUs exhibit certain shortcomings. They lack a certain degree of Well-formedness, and Descriptiveness. Thus, while SGUs and SCUs might appear similar in performance, a closer inspection reveals a slightly better performance of the SCUs.

## 5.6 Discussion

In our research, we found out that there are more effective ways of approximating SCUs than only with STUs. There are several aspects worth discussion. Firstly, it appears that the automatic intrinsic evaluation metric, based on ROUGE-1-F1, exhibits a weak correlation with human evaluation. This raises concerns about the effectiveness of using this metric in previous studies to evaluate the quality of SCU approximations. Secondly, it seems that we may not need the costly SCUs and their approximations to compare summarization systems or rank long summaries (TAC09). Surprisingly, a simple sentence splitting baseline already achieves competitive results compared to SCUs. Finally, SCUs and their approximations offer the most value to rank short summaries (PyrXSum and RealSumm).

## 6 Conclusions

This work primarily focuses on automating the pyramid method. We propose two new methods to approximate SCUs and systematically evaluate the intrinsic quality of SCUs and their approximations. Our experiments on extrinsic evaluation suggests that there might be no need for costly SCUs and their approximations when comparing summarization systems.



## 308 Limitations

## 309 Limitations

310 First, we would have liked to achieve better per-  
311 formance with SMUs generated from AMR, since,  
312 in theory, AMR graph splitting seems ideal to de-  
313 compose a textual meaning into parts, and AMR  
314 generation systems promise to phrase any such sub-  
315 graph in natural language. Inspecting all three parts  
316 of the pipeline (parsing, splitting, and generating),  
317 we find that most issues are likely due to our manu-  
318 ally designed splitting strategy. While the rules are  
319 simple and their creation has profited from com-  
320 munication with AMR-knowledgeable researchers,  
321 the main problem is that there are countless possi-  
322 bilities of how to split an AMR, and the importance  
323 of rules may strongly depend on the further graph  
324 context. Therefore, we believe it is likely that fu-  
325 ture work can strongly improve the AMR approach  
326 by better learning how to better split meaning rep-  
327 resentation graphs.

328 Second, we used an NLI system that was fine-  
329 tuned on gold SCUs extracted from the develop-  
330 ment data (TAC08), since this NLI system was  
331 found to work best by [Zhang and Bansal \(2021\)](#).  
332 While in principle this does not affect the evalua-  
333 tion of SxUs, which was the focus of this paper, it  
334 is not unlikely that by training the NLI system on  
335 each SxU type separately, the results of SxUs may  
336 further improve and so the results for human SCUs  
337 can be considered as slightly optimistic. In general,  
338 the interaction of automatic NLI and SCUs in an  
339 automated pyramid needs to be better understood.  
340 Other recent findings ([Chen and Eger, 2022](#); [Steen  
341 et al., 2023](#)) suggest that NLI models may play  
342 an underestimated role in NLG evaluation. As a  
343 check, we repeated evaluation with an NLI system  
344 without SCU fine-tuning, and observe significant  
345 performance drops across the board, indicating that  
346 i) SCU results are likely not too over-optimistic  
347 in comparison to SxUs; and ii) the effective adap-  
348 tation strategy of the NLI system indeed may be  
349 the second cornerstone of an accurate automatic  
350 pyramid and therefore should be better explored in  
351 future work.

## 352 References

353 Laura Banarescu, Claire Bonial, Shu Cai, Madalina  
354 Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin  
355 Knight, Philipp Koehn, Martha Palmer, and Nathan  
356 Schneider. 2013. Abstract meaning representation

for sembanking. In *Proceedings of the 7th linguis-  
357 tic annotation workshop and interoperability with  
358 discourse*, pages 178–186. 359

Manik Bhandari, Pranav Narayan Gour, Atabak Ash-  
360 faq, Pengfei Liu, and Graham Neubig. 2020. [Re-  
361 evaluating evaluation in text summarization](#). In  
362 *Proceedings of the 2020 Conference on Empirical  
363 Methods in Natural Language Processing (EMNLP)*,  
364 pages 9347–9359, Online. Association for Computa-  
365 tional Linguistics. 366

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
367 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
368 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
369 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
370 Gretchen Krueger, Tom Henighan, Rewon Child,  
371 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
372 Clemens Winter, Christopher Hesse, Mark Chen, Eric  
373 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,  
374 Jack Clark, Christopher Berner, Sam McCandlish,  
375 Alec Radford, Ilya Sutskever, and Dario Amodei.  
376 2020. [Language models are few-shot learners](#). 377

Yanran Chen and Steffen Eger. 2022. Menli: Robust  
378 evaluation metrics from natural language inference.  
379 *arXiv preprint arXiv:2208.07316*. 380

Sebastian Gehrmann, Elizabeth Clark, and Thibault Sel-  
381 lam. 2022. [Repairing the cracked foundation: A sur-  
382 vey of obstacles in evaluation practices for generated  
383 text](#). 384

Karl Moritz Hermann, Tomas Kocisky, Edward Grefen-  
385 stette, Lasse Espeholt, Will Kay, Mustafa Suleyman,  
386 and Phil Blunsom. 2015. Teaching machines to read  
387 and comprehend. In *NIPS*. 388

David M. Howcroft, Anya Belz, Miruna-Adriana  
389 Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad  
390 Mahamood, Simon Mille, Emiel van Miltenburg,  
391 Sashank Santhanam, and Verena Rieser. 2020.  
392 [Twenty years of confusion in human evaluation: NLG  
393 needs evaluation sheets and standardised definitions](#).  
394 In *Proceedings of the 13th International Conference  
395 on Natural Language Generation*, pages 169–182,  
396 Dublin, Ireland. Association for Computational Lin-  
397 guistics. 398

Shashi Narayan, Shay B. Cohen, and Mirella Lapata.  
399 2018. [Don’t give me the details, just the summary!  
400 topic-aware convolutional neural networks for ex-  
401 treme summarization](#). In *Proceedings of the 2018  
402 Conference on Empirical Methods in Natural Lan-  
403 guage Processing*, pages 1797–1807, Brussels, Bel-  
404 gium. Association for Computational Linguistics. 405

Ani Nenkova and Rebecca Passonneau. 2004. [Evaluat-  
406 ing content selection in summarization: The pyramid  
407 method](#). In *Proceedings of the Human Language  
408 Technology Conference of the North American Chap-  
409 ter of the Association for Computational Linguistics:  
410 HLT-NAACL 2004*, pages 145–152, Boston, Mas-  
411 sachusetts, USA. Association for Computational Lin-  
412 guistics. 413

414 Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,  
 415 Jason Weston, and Douwe Kiela. 2020. [Adversarial](#)  
 416 [NLI: A new benchmark for natural language under-](#)  
 417 [standing](#). In *Proceedings of the 58th Annual Meet-*  
 418 *ing of the Association for Computational Linguistics*,  
 419 pages 4885–4901, Online. Association for Computa-  
 420 tional Linguistics.

421 NIST. 2008. [https://tac.nist.gov/](https://tac.nist.gov/publications/2008/papers.html)  
 422 [publications/2008/papers.html](https://tac.nist.gov/publications/2008/papers.html).

423 NIST. 2009. [https://tac.nist.gov/](https://tac.nist.gov/publications/2009/papers.html)  
 424 [publications/2009/papers.html](https://tac.nist.gov/publications/2009/papers.html).

425 OpenAI. 2023. [Gpt-4 technical report](#).

426 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,  
 427 Carroll Wainwright, Pamela Mishkin, Chong Zhang,  
 428 Sandhini Agarwal, Katarina Slama, Alex Ray, et al.  
 429 2022. Training language models to follow instruc-  
 430 tions with human feedback. *Advances in Neural*  
 431 *Information Processing Systems*, 35:27730–27744.

432 Maxime Peyrard, Teresa Botschen, and Iryna Gurevych.  
 433 2017. [Learning to score system summaries for bet-](#)  
 434 [ter content selection evaluation](#). In *Proceedings of*  
 435 *the Workshop on New Frontiers in Summarization*,  
 436 pages 74–84, Copenhagen, Denmark. Association for  
 437 Computational Linguistics.

438 Ori Shapira, David Gabay, Yang Gao, Hadar Ronen, Ra-  
 439 makanth Pasunuru, Mohit Bansal, Yael Amsterdamer,  
 440 and Ido Dagan. 2019. [Crowdsourcing lightweight](#)  
 441 [pyramids for manual summary evaluation](#). In *Pro-*  
 442 *ceedings of the 2019 Conference of the North Amer-*  
 443 *ican Chapter of the Association for Computational*  
 444 *Linguistics: Human Language Technologies, Volume*  
 445 *1 (Long and Short Papers)*, pages 682–687, Min-  
 446 neapolis, Minnesota. Association for Computational  
 447 Linguistics.

448 Julius Steen, Juri Opitz, Anette Frank, and Katja Mark-  
 449 ert. 2023. With a little push, nli models can robustly  
 450 and efficiently predict faithfulness. *arXiv preprint*  
 451 *arXiv:2305.16819*.

452 Shiyue Zhang and Mohit Bansal. 2021. [Finding a bal-](#)  
 453 [anced degree of automation for summary evaluation](#).  
 454 In *Proceedings of the 2021 Conference on Empirical*  
 455 *Methods in Natural Language Processing*, pages  
 456 6617–6632, Online and Punta Cana, Dominican Re-  
 457 public. Association for Computational Linguistics.

## 458 A Appendix

### 459 A.1 Obtaining AMR Sub-graphs

460 Given an AMR graph, we first extract all predicates  
 461 to discern their semantic meaning as we view them  
 462 to form the core of a sentence’s meaning. Subse-  
 463 quently, the argument connections within the pred-  
 464 icates were examined. If a predicate is connected  
 465 to at least one core role (CR), indicated by ARG<sub>n</sub>  
 466 edge label, we extract a sub-graph for every CR of

467 this predicate containing the CR and the underly-  
 468 ing connections. Below we show an example of  
 469 extracting sub-graphs (Figure 3 and Figure 4) from  
 470 an AMR graph (Figure 2) and the corresponding  
 471 SMUs.

472 **Input sentence:** Godfrey Elfwick recruited via  
 473 Twitter to appear on World Have Your Say.

474 **SMUs from sub-graphs:** Godfrey Elfwick was  
 475 recruited. # Godfrey Elfwick will appear on World  
 476 Have Your Say.

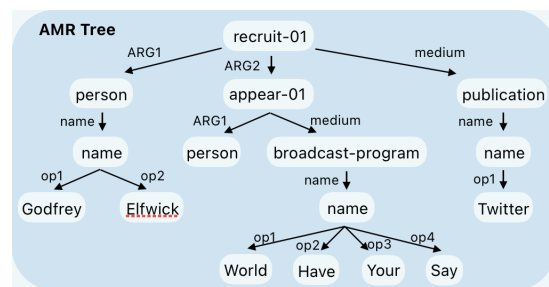


Figure 2: The AMR graph for the sentence “Godfrey Elfwick recruited via Twitter to appear on World Have Your Say”

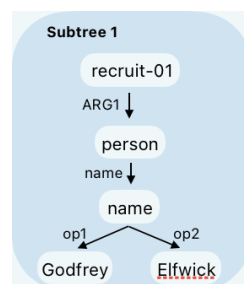


Figure 3: Sub-graph 1

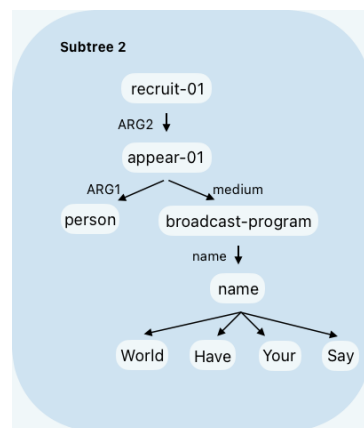


Figure 4: Sub-graph 2

## 477 A.2 Obtaining SGUs from the GPT Models

478 Below we show an example of the prompt we use  
479 to obtain SGUs from the GPT-3.5-Turbo and GPT-  
480 4 models through OpenAI APIs<sup>4</sup>. We randomly  
481 choose a summary and the corresponding human-  
482 written SCUs from the TAC08 dataset as the one-  
483 shot example. We did not alter the standard param-  
484 eters except for temperature which we set to 0 to  
485 ensure reproducibility.

486 **Instruction** You split the provided input in small  
487 sentences separated by an #. The split sentences  
488 represent subsentences of the original sentences.

489 **Example Input** Irish PM Ahern said the main  
490 goal of the US-brokered Good Friday pact of 1998,  
491 a joint Catholic-Protestant administration in North-  
492 ern Ireland, could be revived only with a complete  
493 end of IRA weapons use. The landmark peace  
494 deal led to a virtual end of violence in that area.  
495 Sinn Fein leader Gerry Adams has appealed to IRA  
496 members to end their armed struggle in favor of  
497 democratic politics. Hopes are rising in Northern  
498 Ireland that the IRA will disarm. British PM Blair  
499 and Ahern will chair a review of the Northern Ire-  
500 land situation in London.

501 **Example Output** Good Friday pact was agreed  
502 in 1998 # Good Friday pact was a peace pact #  
503 Good Friday pact set up a joint Catholic-Protestant  
504 administration in Northern Ireland # Good Friday  
505 pact was mediated by the US # Irish Republican  
506 Army increased activity # Irish PM Ahern called to  
507 end violence # Sinn Fein Adams called to end vio-  
508 lence # Hope in Northern Ireland that the IRA will  
509 disarm # British PM Blair and Ahern will chair a  
510 review of the Northern Ireland situation in London

511 **Testing Reference Summary** Netherlands mid-  
512 fielder Wesley Sneijder has joined French Ligue 1  
513 side Nice on a free transfer.

514 **Output From GPT-3.5-Turbo** Netherlands mid-  
515 fielder Wesley Sneijder has joined Nice # Sneijder  
516 was a free transfer # Nice is a French Ligue 1 side  
517 **Output From GPT-4** Netherlands midfielder Wes-  
518 ley Sneijder # Sneijder joined French Ligue 1 side  
519 Nice # Joined on a free transfer

## 520 A.3 Extrinsic Evaluation Details

521 **System-level** correlation assesses the metric’s abil-  
522 ity to compare different summarization systems.  
523 This is denoted as K and measures the correlation

<sup>4</sup><https://openai.com/blog/openai-api>

524 between human scores (h), the metric (m), and the  
525 generated summaries (s) for N examples across  
526 S systems in the meta-evaluation dataset. The  
system-level correlation is then defined as:

$$K_{m,h}^{sys} = K\left(\left[\frac{1}{N} \sum_{i=1}^N m(s_{i1}), \dots, \frac{1}{N} \sum_{i=1}^N m(s_{iS})\right], \left[\frac{1}{N} \sum_{i=1}^N h(s_{i1}), \dots, \frac{1}{N} \sum_{i=1}^N h(s_{iS})\right]\right)$$

527 **Summary-level** correlation assesses the metric’s  
528 ability to compare summaries produced by different  
529 systems for a common document(s). The summary-  
530 level correlation is then defined as:  
531  
532

$$K_{m,h}^{sum} = \frac{1}{N} \sum_{i=1}^N K\left([m(s_{i1}), \dots, m(s_{iS})], [h(s_{i1}), \dots, h(s_{iS})]\right)$$

## 533 A.4 Dataset Details

534 The TAC08 dataset includes 96 examples and out-  
535 puts from 58 systems, while TAC09 contains 88 ex-  
536 amples and outputs from 55 systems. Both datasets  
537 contain multiple reference summaries for each ex-  
538 ample, as well as the corresponding SCU annota-  
539 tions.  
540

541 The REALSumm dataset contains 100 test ex-  
542 amples from the CNN/DM dataset (Hermann et al.,  
543 2015) and 25 system outputs. The SCUs are la-  
544 beled by the authors and SCU-presence labels are  
545 collected using Amazon Mechanical Turk (AMT).

546 PyrXSum (Zhang and Bansal, 2021) includes  
547 100 test examples from the XSum dataset (Narayan  
548 et al., 2018), which contains short and abstractive  
549 summaries. Similar to REALSumm, the SCUs are  
550 manually labeled by the authors and SCU-presence  
551 labels are collected for summaries generated by 10  
552 systems through AMT.  
553

## 554 A.5 Human annotated evaluation

555 The text units of each example were analyzed re-  
556 garding Well-formedness, Descriptiveness and Hal-  
557 lucination. For each dimension, we classified it

558 into one of three categories based on the evalua-  
 559 tor’s satisfaction with the system’s output. These  
 560 categories ranged from “1 - Unhappy with system  
 561 output”, “2 - implying dissatisfaction or a less than  
 562 satisfactory result”, to “3 - Okay with system out-  
 563 put (3)”. Below we denote ASCU for approximated  
 564 summary content unit (e.g., SGUs\_4, SMUs, STUs  
 565 or SCUs) and provide a detail definition for each  
 566 evaluation category:

- 567 • Well-formedness (surface quality)
  - 568 – 1: Many ASCUs are are not concise En-  
569 glish sentences
  - 570 – 2: Some ASCUs are not concise English  
571 sentences
  - 572 – 3: Almost all or all ASCUs are concise  
573 English sentences
- 574 • Descriptiveness (meaning quality I)
  - 575 – 1: Many meaning facts of the summary  
576 have not been captured well by the AS-  
577 CUs
  - 578 – 2: Some meaning facts of the summary  
579 have not been captured by the ASCUs
  - 580 – 3: Almost every or every meaning fact  
581 of the summary has been captured by a  
582 ASCU
- 583 • Hallucination (meaning quality II)
  - 584 – 1: Many ASCUs describe meaning that  
585 is not grounded in the summary
  - 586 – 2: There is some amount of ASCUs that  
587 describes meaning that is not grounded  
588 in the summary
  - 589 – 3: Almost no or no ASCU describes  
590 meaning that is not grounded in the sum-  
591 mary

592 In the following we show two examples of the  
 593 reference summaries and the corresponding AS-  
 594 CUs from PyrXSum and RealSumm, respectively:

- 595 • **Reference summary:** West Ham say they are  
 596 “disappointed” with a ruling that the terms of  
 597 their rental of the Olympic Stadium from next  
 598 season should be made public.
- 599 • **SCUs:** West Ham are “disappointed” with a  
 600 ruling # The ruling is that their rental terms  
 601 should be made public # West Ham will rent  
 602 the Olympic Stadium from next season

- **SMUs:** West Ham say they are disappointed 603  
 by the ruling that their terms of rental for the 604  
 Olympic Stadium next season should be pub- 605  
 lic # The ruling that the terms of West Ham’s 606  
 Olympic Stadium rental next season should 607  
 be public was disappointing # West Ham rent 608  
 the Olympic Stadium # West Ham will rent 609  
 the Olympic Stadium next season 610
- **SGUs\_4:** West Ham is disappointed with 611  
 a ruling # Terms of their Olympic Stadium 612  
 rental should be made public # Olympic Sta- 613  
 dium rental starts next season 614
- **STUs:** West Ham say they are “disappointed” 615  
 with a ruling that the terms of their rental of 616  
 the Olympic Stadium from next season should 617  
 be made public # They are “disappointed” 618  
 with a ruling that the terms of their rental of 619  
 the Olympic Stadium from next season should 620  
 be made public # should made public 621