

Overview of BioASQ 2025: The thirteenth BioASQ challenge on large-scale biomedical semantic indexing and question answering

Anastasios Nentidis¹, Georgios Katsimpras¹, Anastasia Krithara¹, Martin Krallinger², Miguel Rodríguez-Ortega², Eduard Rodríguez-López², Natalia Loukachevitch³, Andrey Sakhovskiy^{5,6}, Elena Tutubalina^{4,5}, Dimitris Dimitriadis⁷, Grigorios Tsoumakas^{7,10}, George Giannakoulas⁷, Alexandra Bekiaridou⁸, Athanasios Samaras⁷, Giorgio Maria Di Nunzio⁹, Nicola Ferro⁹, Stefano Marchesin⁹, Marco Martinelli⁹, Gianmaria Silvello⁹, and Georgios Paliouras¹

¹ National Center for Scientific Research “Demokritos”, Athens, Greece
{tasosnent, gkatsibras, akrithara, paliourg}@iit.demokritos.gr

² Barcelona Supercomputing Center, Barcelona, Spain
{martin.krallinger, mirodrig8, eduard.rodriguez}@bsc.es

³ Moscow State University, Russia
louk_nat@mail.ru

⁴ Artificial Intelligence Research Institute, Russia

⁵ Kazan Federal University, Russia

⁶ SberAI & Skoltech, Russia

{andrey.sakhovskiy, tutubalinaev}@gmail.com

⁷ Aristotle University of Thessaloniki, Greece
{dndimitri,greg}@csd.auth.gr, {g.giannakoulas, th.samaras.as}@gmail.com

⁸ Northwell Health, USA

ampekariidou@gmail.com

⁹ University of Padua, Italy

{name.surname}@unipd.it

¹⁰ Archimedes, Athena Research Center, Greece

Abstract. This is an overview of the thirteenth edition of the BioASQ challenge in the context of the Conference and Labs of the Evaluation Forum (CLEF) 2025. BioASQ is a series of international challenges promoting advances in large-scale biomedical semantic indexing and question answering. This year, BioASQ consisted of new editions of the two established tasks, b and Synergy, and four new tasks: a) *Task Multi-ClinSum* on multilingual clinical summarization. b) *Task BioNNE-L* on nested named entity linking in Russian and English. c) *Task ELCardioCC* on clinical coding in cardiology. d) *Task GutBrainIE* on gut-brain interplay information extraction. In this edition of BioASQ, 83 competing teams participated with more than 1000 distinct submissions in total for the six different shared tasks of the challenge. Similar to previous editions, several participating systems achieved competitive performance, indicating the continuous advancement of the state-of-the-art in the field.

Keywords: Biomedical knowledge · Semantic Indexing · Question Answering

1 Introduction

The BioASQ challenge was introduced over a decade ago, aiming to advance the state-of-the-art in large-scale biomedical semantic indexing and question answering (QA) [77]. To achieve this, it hosts annual shared tasks, creating benchmark datasets that reflect the real-world information needs of biomedical experts. These include new versions of established tasks that remain relevant and timely, as well as novel tasks introduced to explore and address unmet biomedical information needs. These shared tasks provide research teams worldwide, who are developing systems for biomedical semantic indexing and QA, with access to publicly available datasets, a standardized evaluation framework, and opportunities for knowledge exchange through the BioASQ challenge and workshop.

Here, we present the shared tasks and the datasets of the thirteenth edition of the BioASQ challenge in 2025, as well as a condensed overview of the participating systems and their performance. The remainder of this paper is organized as follows. First, Section 2 presents a general description of the shared tasks, which took place in 2025, and the corresponding datasets developed for the challenge. Then, Section 3 provides a brief overview of the participating systems for the different tasks. Detailed descriptions for some of the systems are available in the respective extended overviews of each task and the proceedings of the BioASQ lab. Subsequently, in Section 4, we present the performance of the systems for each task, based on state-of-the-art evaluation measures or manual assessment. Finally, in Section 5 we draw some conclusions.

2 Overview of the tasks

The thirteenth edition of the BioASQ challenge consisted of six tasks [50]: (i) *Task b* on biomedical semantic question answering. (ii) *Task Synergy* on question answering developing biomedical topics. (iii) *Task MultiClinSum* on multilingual clinical summarization. (iv) *Task BioNNE-L* on nested named entity linking in Russian and English. (v) *Task ELCardioCC* on clinical coding in cardiology. (vi) *Task GutBrainIE* on gut-brain interplay information extraction. In this section, we first describe this year’s editions of the two established tasks b (task 13b) and Synergy (Synergy 13) [54] with a focus on differences from previous editions of the challenge [51, 57]. Additionally, we also introduce the four new BioASQ tasks, MultiClinSum [66], BioNNE-L [67], ELCardioCC [18], and GutBrainIE [48].

2.1 Task 13b

BioASQ *task 13b* is the thirteenth edition of the established BioASQ *task b* on Biomedical QA [55]. This year, it took place in three phases: i) Phase A:

biomedical questions in English were provided, and the systems had to retrieve relevant material (PubMed documents and snippets). ii) Phase A+, the systems had to provide ‘exact’ and ‘ideal’ answers. Depending on question type, the ‘exact’ answer can be a *yes* or *no* (yes/no), an entity name, such as a disease or gene (factoid), or a list of entity names (list). The ‘ideal’ answer is a paragraph-sized summary, regardless of question type. iii) Phase B: Some relevant material was provided for each question, selected by the BioASQ experts, and the systems had to provide new answers given this additional information.

About 340 new biomedical questions annotated with golden documents, snippets, and answers (‘exact’ and ‘ideal’), were developed for testing. In addition, a training set of 5,389 biomedical questions, accompanied by answers, and supporting evidence (documents and snippets), was available from previous versions of the tasks, as a unique resource for the development of question-answering systems [34]. Table 1 presents some statistics of both training and test datasets for task 13b. The test data for task 13b were split into four independent bi-weekly batches consisting of 85 questions each, as presented in Table 1.

Table 1. Statistics on the training and test datasets of task 13b. The numbers for the documents and snippets refer to averages per question.

Batch	Size	Yes/No	List	Factoid	Summary	Documents	Snippets
Train	5389	1459	1047	1600	1283	9.74	12.78
Test 1	85	17	23	26	19	2.68	3.74
Test 2	85	17	19	27	22	2.71	3.06
Test 3	85	22	22	20	21	3.00	3.66
Test 4	85	26	19	22	18	3.15	3.92
Total	5729	1541	1130	1695	1363	9.33	12.23

2.2 Task Synergy 13

BioASQ *task Synergy* was originally introduced in 2020 with the aim of promoting research in developing biomedical topics, such as COVID-19 [36, 37]. The design of this task as an ongoing dialogue allows experts to pose open-ended questions for developing topics, for which they do not know in advance whether a definitive answer can be given, in order to obtain relevant material (documents and excerpts) retrieved from the systems. After assessing this material, they provide feedback to the systems on its relevance and on whether it is sufficient to answer their question, by marking respective questions as *ready to answer*. This process is repeated iteratively in rounds with new material considered in each round, based on updates to the original document resource [56]¹¹. For *ready to answer* questions, they receive exact and ideal answers as well, assess them, and provide feedback that can be used by the systems to improve their responses to

¹¹ As of 2023, this evolving document resource is PubMed [51]

these questions in the remaining rounds. The experts can also mark a question as *closed* if they receive a fully satisfactory answer that is not expected to change or if they are no longer interested in the question.

A training dataset of 366 questions on developing topics with incremental annotations with relevant material and answers is already available from previous versions of *task Synergy* [54, 53, 59, 58]. During the *task Synergy 13*, this set was extended with 47 new questions on developing health topics, such as infectious, rare, and genetic diseases, and women’s and reproductive health. Meanwhile, 27 questions from the previous version of the task remained open and were enriched with more recent evidence and updated answers [55]. Overall, 74 questions were considered in the four rounds of *task Synergy 13*. The number of yesno, list, factoid, and summary questions was 23, 19, 14, and 18, respectively.

2.3 Task MultiClinSum

There is a rapid accumulation of various types of clinical content, including medical records and publications such as clinical case reports, written not only in English but also in many other languages. Some clinical reports can be very lengthy, making it challenging for healthcare professionals and even patients to comprehend and extract key clinical insights. Large Language Models (LLMs) have shown promising results in automatic summarization, helping to condense lengthy clinical documents into shorter versions or summaries that retain the most relevant clinical information. Therefore, there is a pressing need to evaluate and benchmark the performance of different clinical summarization methods, especially for content written in multiple languages.

We introduce the *MultiClinSum* task covering the automatic summarization of lengthy clinical case reports written in English, Spanish, French, and Portuguese. The *MultiClinSum* task relies on a corpus of manually selected full clinical case reports with their corresponding summaries derived from case report publications written in the mentioned languages. This gold standard dataset comprises 1,280 pairs of full-text and summary in English, 534 in Spanish, 200 in Portuguese and 200 in French. To increase the size of the dataset, both full-text and summary texts of each language were translated with neural machine translation models into the other languages, resulting in a total of 1,976 pairs for each language. An additional large-scale dataset was also created derived from the PMC-Patients clinical cases (full-text) and their corresponding summary extracted from the PubMed abstracts. Table 2 shows the corpus statistics for each sub-track of MultiClinSum.

For the evaluation assessment, the automatically generated summaries were compared with the summaries that had been manually generated by the original authors, using Rouge-L [41] scores and BERTScore [85]. As clinical case reports do share commonalities with medical discharge summaries (patient demographics, relevant medical history, clinical presentation, diagnostic process, intervention, treatment, outcome and follow-up), insights provided by the *MultiClinSum* results can be of practical relevance also for clinical records summarization scenarios.

Table 2. Statistics for the datasets provided for MultiClinSum indicating the number of fulltext-summary pairs available for each sub-track.

Sub-track	Lang.	Dataset	Fulltext-summary pairs
MultiClinSum-gs-en	EN	Native/Transl. gold stand.	988
MultiClinSum-gs-es	ES	Native/Transl. gold stand.	988
MultiClinSum-gs-fr	FR	Native/Transl. gold stand.	1061
MultiClinSum-gs-pt	PT	Native/Transl. gold stand.	1034
MultiClinSum-ls-es	EN	Large Scale (PMC-Patients)	28.902
MultiClinSum-ls-es	ES	Large Scale (PMC-Patients)	28.902
MultiClinSum-ls-fr	FR	Large Scale (PMC-Patients)	28.902
MultiClinSum-ls-pt	PT	Large Scale (PMC-Patients)	28.902

2.4 Task BioNNE-L

In the BioNNE-L Shared Task [67], we address the medical entity linking task, also known as Medical Concept Normalization (MCN), which is to map given entities to the most relevant vocabular entries from an external source, e.g., concepts from the UMLS metathesaurus [7] identified with concept unique identifiers (CUIs). Although the task has been widely explored in recent years, existing approaches usually treat each entity individually, medical entities often form a nested structure, where an entity can be a subpart of another entity. One of the key features of BioNNE-L is the focus on nested entities that are (i) derived from the MCN annotation of the NEREL-BIO corpus [45, 46] and (ii) supplemented by newly annotated data in both English and Russian. The annotated entity types are disorders (*DISO*), anatomical structures (*ANAT*), and chemicals (*CHEM*) normalized to UMLS. The competition was organized into three subtasks that fell under two evaluation tracks: 1. **Monolingual track** that treated English and Russian data independently; 2. **Bilingual track** that required a single bilingual model for the combined Russian and English data. Data statistics for both tracks, as well as the normalization dictionary, are summarized in Table 3.

All BioNNE-L materials can be found on the shared task’s GitHub¹² and Codalab pages¹³. Annotated data and normalization dictionary are also available at HuggingFace¹⁴.

2.5 Task ELCardioCC

Cardiovascular diseases affect a significant portion of the global population, accounting for 32% of global deaths according to WHO¹⁵. Automated clinical coding plays a crucial role in transforming unstructured real-world medical data gathered from patients into structured information, in order to facilitate clinical research and analysis. However, existing research predominantly focuses on

¹² <https://github.com/nerel-ds/NEREL-BIO/tree/master/BioNNE-L.Shared.Task>

¹³ <https://codalab.lisn.upsaclay.fr/competitions/21568>

¹⁴ <https://huggingface.co/datasets/andorei/BioNNE-L>

¹⁵ <https://www.who.int/health-topics/cardiovascular-diseases>

Table 3. BioNNE-L 2025 statistics for Disorder (**DISO**), Chemical (**CHEM**), and Anatomical Structure (**ANAT**) among Russian and English entities as well as normalization dictionary statistics.

Entity type	Refined NEREL-BIO				Novel data		Dictionary	
	Train		Dev		Test			
	Ru	En	Ru	En	Ru	En	Ru	En
# documents	716	54	50	50	154	154	—	—
Number of entities								
DISO	11,168	1,200	925	1,029	2,811	3,068	91,867	1,825,048
CHEM	4,741	579	531	564	1,218	1,345	47037	1,732,096
ANAT	8,346	911	878	901	2,186	2,248	6899	345,043
	24,255	2,690	2,334	2,494	6,215	6,661	145,803	3,902,187

English clinical text, leaving other languages, such as Greek, underrepresented. To this end, we propose a new *ELCardioCC* task [18], which concerns i) the assignment of cardiology-related ICD-10 codes to discharge letters from Greek hospitals, ii) the extraction of the specific mentions of ICD-10 codes from the discharge letters.

In detail, the participants in the *ELCardioCC* task were tasked with developing named entity recognition (NER), entity linking (EL) and multi-label classification - explainable AI (MLC-X) systems using a specialized corpus of discharge letters. These discharge letters, which were written in Greek contained valuable medical information about patients’ conditions, treatments, and outcomes. The corpus was meticulously annotated with the positions of mentions (such as chief complaint, diagnosis, prior medical history, drugs and cardiac echo) and their corresponding ICD-10 codes. The training dataset includes 1,000 discharge letters, while the test set comprises 500 letters. System performance was evaluated using the micro F1 score.

2.6 Task GutBrainIE

Recent scientific evidence suggests a connection between *brain-related diseases* and the *gut microbiota* that may play a critical role in mental health-related disorders or diseases like Parkinson’s, and Alzheimer’s [3, 11, 14, 23]. The scientific literature on this topic is rapidly expanding, making it increasingly challenging for clinicians and researchers to stay up to date. For example, in 2020, approximately 200 articles were published on the relationship between gut microbiota and mental health; by 2024, this number had more than doubled to over 450 publications. The *GutBrainIE* Task aims to foster the development of Information Extraction (IE) systems that support experts by automatically extracting and linking knowledge from biomedical abstracts, facilitating the understanding of gut-brain interplay and its role in mental health and neurological diseases.

The *GutBrainIE* Task comprises four subtasks of increasing difficulty: Named Entity Recognition (NER), which identifies and classifies entity mentions in

PubMed abstracts about the gut-brain interplay focusing on mental health and the Parkinson’s disease; Binary Tag-based Relation Extraction (BT-RE), which detects whether pairs of entities are in relation without specifying the relation type; Ternary Tag-based Relation Extraction (TT-RE), which extends BT-RE by also assigning a relation label to each related pair; and Ternary Mention-based Relation Extraction (TM-RE), which further localizes the exact entity mentions involved in each relation and assigns the appropriate relation label.

The dataset includes over 1000 documents with annotated entity mentions and relations, organized into Training, Development, and Test sets. The train set is further divided into four quality tiers: expert-curated (Platinum), expert-annotated (Gold), student-annotated (Silver), and automatically generated (Bronze). Development and Test sets contain only expert annotations (Platinum+Gold).

Table 4. Dataset statistics for *GutBrainIE*.

Collection	# Docs	# Entities	Ents/Doc	# Rels	Rels/Doc
Train Platinum	111	3638	32.77	1455	13.11
Train Gold	208	5192	24.96	1994	9.59
Train Silver	499	15275	30.61	10616	21.27
Train Bronze	749	21357	28.51	8165	11.90
Development Set	40	1117	27.93	623	15.58
Test Set	40	1237	30.92	777	19.42

3 Overview of participation

Overall, 83 distinct teams participated in the thirteenth edition of the BioASQ challenge, submitting more than 1000 distinct runs for the six different shared tasks of the challenge. The majority of the teams focused on a single task, still some of them participated in two or even three BioASQ tasks¹⁶. In this section, we provide a condensed overview of the methods developed by the participating teams for each of the BioASQ tasks. However, a more detailed overview of these methods will be available in the extended overview of each task [55, 66, 67, 18, 48], and some method-specific descriptions will be available in the proceedings of the thirteenth BioASQ workshop¹⁷.

3.1 Task 13b

This year, 46 teams participated in task 13b, submitting a total of 734 different submissions generated by 146 distinct systems across all four batches for the

¹⁶ In particular, two teams participated in 13b & Synergy13, one in 13b & MultiClinSum, one in 13b & GutBrainIE, one in 13b, MultiClinSum, & ElCardioCC, and one in BioNNE-L, ElCardioCC, & GutBrainIE

¹⁷ <https://www.bioasq.org/workshop2025/proceedings>

three phases A, A+, and B. This corresponds to a significant increase in participation, compared to the 26 teams in the previous version of the task (12b)[54], which highlights that the task remains timely and relevant. Specifically, 34, 20, and 26 teams competed in phases A, A+, and B of task 13b, with 95, 79, and 88 distinct systems, respectively. Eleven of these teams were involved in all three phases. As in previous years, the open-source system OAQA [82], which achieved top performance in older editions of BioASQ [35], was used as a baseline for phase B *exact answers*.

The participating teams employed a range of well-established and sophisticated techniques. Many teams utilized traditional document retrieval methods such as BM25 and dense retrieval models (e.g. BGE-M3 and MiniLM), often improving results with re-ranking techniques. Some teams incorporated Retrieval-Augmented Generation (RAG) frameworks, using Large Language Models (LLMs) such as Llama, Gemma, GPT, Claude, and Mistral to generate responses. Beyond these methods, the teams also experimented with self-feedback mechanisms, zero-shot and few-shot prompting, ensemble methods, and the integration of biomedical knowledge bases to improve overall performance [63, 81, 72, 6, 33, 21, 74, 2, 4, 30, 9].

3.2 Task Synergy 13

In the thirteenth edition of BioASQ, five teams participated in the Synergy task (Synergy 13). These teams submitted 46 runs from 21 distinct systems. Two of these teams participated in task 13b as well, while the remaining three focused exclusively on task Synergy 13. The participating teams primarily utilized LLMs, such as DeepSeek-R1 and Llama. To further enhance performance, the teams experimented with RAG frameworks and employed techniques such as optimized prompting, NER, and majority voting to refine their results [19, 63]. More detailed descriptions for some of the systems are available at the proceedings of the workshop.

3.3 Task MultiClinSum

In general, there has been a very satisfactory participation in the task with promising results in each of the presented sub-tracks. 56 teams registered for the MultiClinSUM task, out of which 11 teams submitted at least one run of their predictions. Specifically, 7 teams participated in the English sub-track, 5 teams in the Spanish, 4 teams in French, and 5 in Portuguese. Each team was allowed to submit up to 5 runs per sub-track. As expected, the best results were obtained in the English sub-track (MultiClinSum-en), which had the highest level of participation. Nevertheless, the others sub-tracks were quite well represented in terms of both participation and novel methodologies applied [66].

3.4 Task BioNNE-L

In total, we’ve received 23 Codalab registrations for the BioNNE-L task, with 7 teams submitting predictions during the evaluation phase. The systems submitted by the participants are summarized in Table 5.

Table 5. Overview of the approaches presented by participants for the BioNNE task. EN stands for the English-oriented and RU for the Russian-oriented tracks.

Team	Track	Approach
verbanexialab	EN	SapBERT w/ lexical and semantic reranking
LYX_DMIIP_FDU	Bilingual,EN,RU	BERGAMOT fine-tuning
BlancaPlanca	Bilingual,EN,RU	BERGAMOT w/ language-specific preprocessing
MSM Lab	Bilingual,EN,RU	Two-step retrieval and ranking pipeline
dstepakov	Bilingual,RU	RoBERTa fine-tuning with contrastive learning
ICUE	Bilingual,EN,RU	BERT, BioSyn, LLM 0-shot reranking
NLPIMP	Bilingual	Russian LaBSE model pre-trained on medical data

Team **verbanexialab** [64] leveraged a SapBERT¹⁸ [42], pre-trained on UMLS concepts, to obtain entity embeddings, followed by a multicomponent re-ranking. They combined embedding cosine similarity with Jaccard similarity for lexical overlap recognition and Levenshtein distance for character-level alignment.

Team **LYX_DMIIP_FDU** [44] fine-tuned a BERGAMOT¹⁹ [69] model for each task via contrastive learning using the train- and dev-set entities to enrich the original vocabularies. The textual context of each entity was used as additional input to enhance the entity representation.

Team **BlancaPlanca** [10] used BERGAMOT for zero-shot retrieval based on entity-concept cosine similarity. They apply language-specific lemmatization for Russian and speed up the inference by chunking the normalization dictionary into type-specific parts of 100k entries each.

Team **MSM Lab** [40] adopted SapBERT [42, 43] and BioMedBERT [25] for two-step retrieval and ranking.

Team **dstepakov** performed the nearest-neighbor search based on the cosine similarity of RoBERTa embeddings [87], fine-tuned contrastively on anchor-positive-negative term triplets via the InfoNCE objective [61].

Team **ICUE** [15] fine-tuned BioSyn [73] using the vocabularies reduced to less than 100k entries each. They fine-tune a separate BERT-based model [17] for

¹⁸ <https://huggingface.co/cambridgeltl/SapBERT-from-PubMedBERT-fulltext>

¹⁹ <https://huggingface.co/andorei/BERGAMOT-multilingual-GAT>

English [5], Russian²⁰, and multilingual [76] tracks, respectively. They re-ranked the initial retrieval results using *DeepSeek-R1-Distill-Llama-8B*²¹.

Team **NLPIMP** performed the zero-shot ranking using a Russian LaBSE [20] model²² pre-trained contrastively on an in-house Russian medical corpus.

3.5 Task ELCardioCC

The ELCardioCC task engaged five teams across its subtasks: NER, EL, and MLC-X, with a total of 13-14 systems submitted for each subtask in addition to baseline models. Most participating systems predominantly utilized transformer-based architectures, especially BERT variants and large multilingual language models (LLMs), for all three tasks. Common approaches included fine-tuning models like Greek BERT and XLM-Roberta for NER, employing semantic similarity with embedding models for EL, and using LLMs for classification and justification in MLC-X, often leveraging cross-lingual techniques to process Greek medical texts.

The ELCardioCC baselines, designed for clarity and reproducibility, primarily used multilingual BERT models adapted for each specific task. The NER baseline involved a fine-tuned cased mBERT model with BIOES tagging. For EL, a context-aware hierarchical classifier built on mBERT was used, reflecting the ICD-10 taxonomy. The MLC-X baseline employed a Greek-BERT model for multi-label classification of the 40 most frequent ICD-10 codes, with variations for document-level predictions and rule-based justification of code selections.

3.6 Task GutBrainIE

The *GutBrainIE* task registered 17 teams submitting runs. Among these, 16 teams participated in NER, 12 in BT-RE, 13 in TT-RE, and 13 in TM-RE. Overall, a total of 391 runs were submitted: 101 for NER, 100 for BT-RE, and 95 for both TT-RE and TM-RE.

Most teams adopted supervised fine-tuning or transformer-based models pre-trained on biomedical text for the NER task [1, 16, 27, 32, 38, 44, 49, 62, 65, 75]. Standard backbones included PubMedBERT, BioBERT, BioLinkBERT, and ELECTRA [12, 25, 39, 83]. Specialized NER architectures, such as GLiNER [84], were also utilized and fine-tuned. Many groups trained multiple models with different random seeds to boost robustness and ensembled their outputs. All teams used platinum, gold, and silver collections for training. A few also used the noisier bronze set, employing cleaning or re-weighting approaches and integrating PubMed data augmentation.

Across the RE subtasks, participants primarily used biomedical pre-trained language models fine-tuned on entity-marker augmented inputs [1, 16, 27, 32, 38, 44, 49, 62, 65, 75]. Among these, the most widely employed include: SapBERT,

²⁰ <https://huggingface.co/KoichiYasuoka/bert-base-russian-upos>

²¹ <https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B>

²² <https://huggingface.co/sergeyzh/LaBSE-ru-turbo>

PubMedBERT, BioBERT, RoBERTa, and ELECTRA [42, 25, 39, 87, 12]. Several teams reformulated RE as a seq2seq problem using REBEL-large [29], directly generating relation tuples or tagged spans in a single pass for all three sub-tasks. Few others leveraged model ensembling and trained with negative-pair subsampling to counter class imbalance and increase generalization capabilities. Finally, some groups experimented with few-shot or Retrieval-Augmented Generation (RAG), prompting large language models to extract both entities and relations with minimal fine-tuning [13, 26, 31, 38].

4 Results

4.1 Task 13b

This section presents the evaluation measures and preliminary results for Task 13b. The evaluation in *task 13b* is done manually by the experts that assess each system response and automatically by employing a variety of established evaluation measures [47] as in the previous versions of the task [52]. Table 6 provides a brief overview of the official measures per response and question type. The results reported for *task 13b* are preliminary, as the final results will be available after the manual assessment of all system responses by the BioASQ team of experts and the enrichment of the ground truth with potential additional relevant items (i.e. documents and snippets), answer elements, and/or synonyms, which is still in progress. The online results pages for Phase A²³, Phase A+²⁴, and Phase B²⁵ will be updated with the final results when available.

The overall performance of the participating systems in document and snippet retrieval (Phase A) per batch of task 13b is presented in Table 7. Both the average and the top performance of the systems seem to drop in the last batch, indicating that the questions in this batch are more challenging for the systems. This could be related to the composition of this batch, which included more questions developed by new BioASQ experts, who have not contributed significantly to the development of the training dataset.

The top performance of the participating systems in *exact answer* generation per question type is presented in Figure 1 for both Phase A+ and Phase B of task 13b, in comparison to the respective performance in the previous version (12b). These preliminary 13b results for phase B suggest that the top systems achieved scores comparable or higher to those of 12b in answering all types of questions (solid lines). In Phase A+ (dashed lines), the top performance is lower, as expected; however, for yesno questions in particular, it is very close to those of Phase B, revealing the increased capability of LLM-based models to address these questions even without being provided ground-truth relevant documents and snippets. These results probably underestimate the performance of the top 13b systems in factoid and list questions, as the preliminary ground truth may

²³ <https://participants-area.bioasq.org/results/13b/phaseA/>

²⁴ <https://participants-area.bioasq.org/results/13b/phaseAplus/>

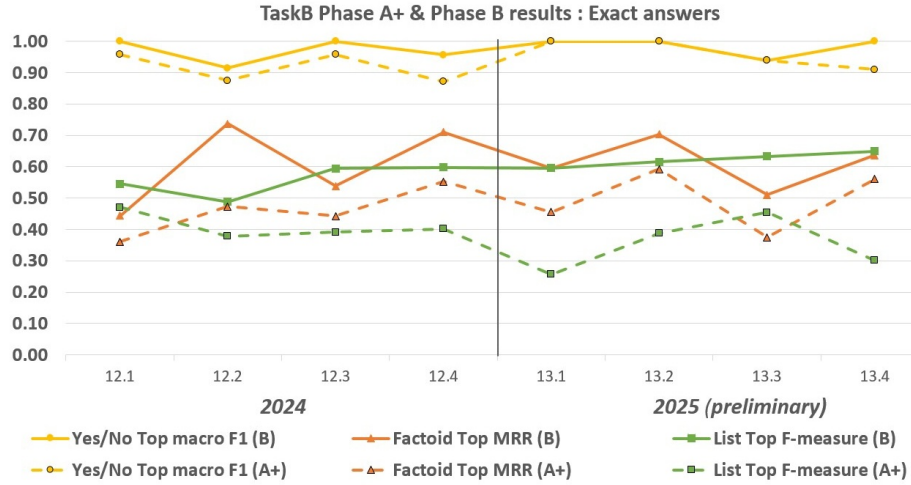
²⁵ <https://participants-area.bioasq.org/results/13b/phaseB/>

Table 6. The evaluation measures for *task 13b* per response type and question type [47].

Resp. type (Phase)	Quest. type	Official measure
Documents (A)	All	Mean Average Precision (MAP)
Snippets (A)	All	F1 (based on character overlaps)
Exact ans. (A+ & B)	List	F1
	Yesno	macro F1 on “yes” & “no” classes
	Factoid	Mean Reciprocal Rank (MRR)
Ideal ans. (A+ & B)	All	Manual scores for precision, recall, repetition, readability

Table 7. The average and top scores of participant systems in Phase A, task 13b.

Batch	Documents		Snippets	
	Average MAP	Top MAP	Average F1	Top F1
1	0.231	0.425	0.053	0.120
2	0.283	0.442	0.084	0.179
3	0.175	0.324	0.052	0.110
4	0.072	0.180	0.024	0.079

**Fig. 1.** The scores of the top systems in *exact answer* generation, for Phase A+ (dashed lines) and B (solid lines), across the test sets of *task 13b* and *task 12b* [60].

miss some synonyms or alternative terms submitted by the participants. Such synonyms will be detected during the enrichment process and will be considered for the final results.

4.2 Task Synergy 13

In *task Synergy 13* we use the same evaluation measures described for *task 13b*, considering only new material for the information retrieval part, an approach known as *residual collection evaluation* [70]. In addition, due to the developing nature of the topics, no answer is available for all of the open questions in each round. Therefore, only the questions indicated as “answer ready” were evaluated for *exact* and *ideal answers* per round.

Table 8 presents the top performance achieved by participating systems per round, in task Synergy 13, for all types of responses. During the four rounds of Synergy 13, the systems managed to identify enough relevant material to provide an answer to 59 of the 74 questions (about 80%). In addition, they also managed to provide at least one *ideal answer* which was considered of ground-truth quality by the respective expert for 35 questions (about 47%). Overall, this dialogue between question-answering systems and biomedical experts allowed the progressive gathering of relevant documents and snippets and the generation of *exact* and *ideal answers* for open questions on developing topics, such as infectious, rare, and genetic diseases, and women’s and reproductive health.

Table 8. The number of “Answer Ready” (AR) questions and the top system performance per round (R) in Task Synergy 13. Retrieval of documents (Top MAP) and snippets (Top F1). Generation of exact factoid (Top MRR), list (Top F1), and yesno (Top ma-F1) answers.

R	AR	Top MAP	Top F1 Snip.	Top MRR	Top F1 list	Top macro-F1
1	19	0.41	0.31	0.67	0.09	1
2	33	0.41	0.29	0.43	0.25	1
3	49	0.46	0.15	0.5	0.26	1
4	55	0.47	0.25	0.45	0.35	1

4.3 Task MultiClinSum

The automatic evaluation of the MultiClinSum results was performed using both BERTScore [85] and ROUGE-LSum metrics. Given the abstractive nature of the task, BERTScore was prioritized as the primary metric due to its superior capacity to capture semantic similarity between generated and reference summaries by leveraging contextualized embeddings, thus effectively recognizing meaning-preserving paraphrases and diverse lexical choices [85].

In contrast, ROUGE-LSum—a sentence-level variant of the ROUGE metric [41]—provides informative measures of summary quality by assessing the overlap of longest common subsequences between candidate and reference summaries. This offers valuable insights into the coverage and faithfulness of the generated content with respect to the original text.

While ROUGE-LSum remains limited by its reliance on surface-level n-gram matching, its inclusion alongside BERTScore ensures a complementary perspective on summary quality, balancing semantic and lexical overlap considerations. The latter, however, was the prioritized metric for submission ranking purposes.

Table 9. Results of the MultiClinSum for each sub-track. Only the top-2 best teams are presented. The best result is in bold.

Team Name	Subtrack	BERTScore			ROUGESum		
		P	R	F1	P	R	F1
seemdog	English	0.8795	0.8608	0.8698	0.3404	0.2398	0.267
pjmath. [79]	English	0.8821	0.8466	0.8637	0.4077	0.2343	0.2805
ggrazhdans [24]	Spanish	0.7699	0.747	0.7578	0.3639	0.2667	0.2899
pjmath. [79]	Spanish	0.7675	0.7392	0.7525	0.3684	0.2703	0.292
pjmath. [79]	French	0.7692	0.7459	0.7567	0.3481	0.2684	0.2843
BU team	French	0.7248	0.7396	0.7315	0.2415	0.289	0.2466
pjmath. [79]	Portuguese	0.7644	0.7377	0.7502	0.35	0.2605	0.2803
ETS-PUCPR [71]	Portuguese	0.7403	0.7351	0.737	0.2802	0.250	0.249

BERTScore results in Table 9 for non-english languages is impaired by the fact that a multilingual bert model is used instead of a language specific encoder model. Participants were able to submit up to 5 runs, the best of which was selected for each subtrack. There was a total of 15, 14, 9 and 7 runs for English, Spanish, French and Portuguese respectively.

4.4 Task BioNNE-L

Following prior research on entity linking [46, 42, 68, 69], we address BioNNE-L as a retrieval task: given a mention, a model must retrieve the top-k concepts from the given UMLS dictionary and employ two ranking-based evaluation metrics: (i) Accuracy@k and (ii) Mean Reciprocal Rank (MRR). Accuracy@k: Accuracy@k=1 if the correct UMLS CUI is retrieved at rank $\leq k$, and Accuracy@k=0 otherwise. $MRR = \frac{1}{|E|} \sum_{e \in E} \frac{1}{rank_e}$, where E is the set of entities, $|E|$ is the number of entities, $rank_e$ is the rank of entity e 's the first correctly retrieved concept among the top k retrieved concepts. As baseline, we adopt zero-shot ranking using BERGMOT [69] with each entity type processed independently to reduce memory footprint caused by extensive dictionary.

The official evaluation results, ordered by Accuracy@1 value, for BioNNE-L are summarized in Table 10. Most of the participants adopted various BERT-based [17] with the top-performance achieved by domain-specific models, such as BERGMOT, SapBERT. Specifically, Team LYX_DMHIP_FDU ranked the first in the multilingual track and the second rank in the two monolingual track by fine-tuning BERGMOT. Top 1 results for the Russian and English data are achieved by multilingual BERGMOT (Team BlancaPlanca) and English SapBERT (Team verbanexialab) models, respectively.

Table 10. Official evaluation results of the BioNNE-L task for the multilingual and monolingual tracks in terms of Accuracy@1 (@1), Accuracy@5 (@5), and MRR. The best results for each track and metric are highlighted in **bold**.

Team	Multilingual				English				Russian			
	#	@1	@5	MRR	#	@1	@5	MRR	#	@1	@5	MRR
verbanexialab	—	—	—	—	1	0.70	0.80	0.74	—	—	—	—
LYX_DMIIP_FDU	1	0.68	0.84	0.75	2	0.66	0.84	0.74	2	0.71	0.84	0.76
BlancaPlanca	2	0.67	0.81	0.73	3	0.64	0.83	0.72	1	0.72	0.83	0.76
MSM Lab	3	0.63	0.76	0.69	4	0.64	0.82	0.71	4	0.65	0.74	0.69
dstepakov	4	0.63	0.71	0.66	—	—	—	—	3	0.70	0.76	0.72
ICUE	5	0.58	0.76	0.66	6	0.51	0.79	0.62	5	0.62	0.72	0.67
baseline	6	0.53	0.70	0.60	5	0.57	0.78	0.66	6	0.52	0.59	0.55
NLPIMP	7	0.41	0.58	0.48	—	—	—	—	—	—	—	—

4.5 Task ELCardioCC

The results of the participants for the ELCardioCC task are presented in tables 11, 12 and 13. For each subtask, the table displays one system per team that achieved the highest F1-score.

Table 11. Performance of participating systems in the ELCardioCC Named Entity Recognition (NER) subtask. Results are reported using micro-averaged precision, recall, and F1 score.

Team	System	Recall	Precision	Micro-F1
bhuang	5nm	0.6448	0.5205	0.5761
droidlyx	system1	0.7059	0.7618	0.7328
ELCardioCC_baseline	mbert_baseline	0.6959	0.7460	0.7201
pjmathematician	config1	0.2484	0.2586	0.2534
enigma	greek-bert-exact-bge-m3	0.7012	0.7328	0.7167

Table 12. Performance of participating systems in the ELCardioCC Entity Linking (EL) subtask. Results are reported using micro-averaged precision, recall, and F1 score.

Team	System	Recall	Precision	Micro-F1
bhuang	5nm	0.5927	0.4852	0.5336
droidlyx	system1	0.6529	0.7046	0.6778
ELCardioCC_baseline	EL_baseline	0.6476	0.6942	0.6701
pjmathematician	config1	0.0616	0.0642	0.0629
enigma	greek-bert-exact-bge-m3	0.6548	0.6844	0.6693

The team droidlyx [44] consistently achieved the highest micro-F1 scores across all subtasks, leading NER (0.733), EL (0.678), and performing strongly

Table 13. Performance of participating systems in ELCardioCC Subtask 3a (Multi-label Classification) and Subtask 3b (Explainable AI). Metrics include Precision (P), Recall (R), and Micro-F1 score. A dash (-) indicates that the system did not participate in the corresponding subtask.

Team	System	Subtask 3a (MLC)			Subtask 3b (X)		
		P	R	F1	P	R	F1
ELCardioCC_baseline	MLCX1_baseline	0.9339	0.7422	0.8271	-	-	-
	MLCX2_baseline	0.9531	0.5864	0.7261	0.6050	0.4442	0.5122
bhuang	1nm	0.6205	0.7676	0.6863	-	-	-
droidlyx	system1	0.8569	0.8377	0.8472	-	-	-
kbogas	w2l_cb	0.2115	0.3421	0.2614	-	-	-
pjmathematician	config4	0.6056	0.2257	0.3288	0.2326	0.0932	0.1331
	config5	0.5860	0.2656	0.3655	-	-	-

in MLC-X (0.847). The ELCardioCC baseline remained highly competitive, particularly excelling in MLC-X with a micro-F1 of 0.827 and maintaining strong performance in NER and EL. Enigma’s systems [80] were consistently near the top in NER and EL but did not participate in MLC-X, while bhuang [28] showed promising results especially in classification, though with more variability. Finally, the pjmathematician [78] system demonstrated significantly low performance .

4.6 Task GutBrainIE

Submitted runs were evaluated using micro- and macro-averaged precision, recall, and F1-score, with micro-F1 used as the reference measure for the leaderboards since it is better suited when classes are imbalanced.

We adopted a baseline employing a fine-tuned NuNER model for NER [8] and a fine-tuned ATLOP model for all RE subtasks [86]. The baseline has also been used to annotate the bronze collection automatically.

Tables 14-17 show each team’s top run beating the baseline for NER, TB-RE, TT-RE, and TM-RE, respectively. The performance gap between the subtasks is noticeable. While NER obtained a top micro-F1 of 0.84 utilizing pretrained biomedical transformers, RE subtasks were more challenging: both TB-RE and TT-RE peaked at approximately 0.65-0.69 micro-F1, and TM-RE reached only 0.46 micro-F1, demonstrating the added difficulty of simultaneously locating and labeling entities and identifying relations among these.

5 Conclusions

This paper provides an overview of the thirteenth BioASQ challenge. This year, BioASQ consisted of six tasks: (i) *Task 13b* on biomedical semantic question answering. (ii) *Task Synergy13* on question answering for developing biomedical topics. (iii) *Task MultiClinSum* on multilingual clinical summarization. (iv) *Task*

Table 14. Performance metrics of each team’s top run beating the baseline for NER. The best result is in bold, the second-best is underlined (micro-averaged).

Team ID	Run name	Precision	Recall	F1
GutUZH [27]	AugEnsemble	0.8384	0.8432	0.8408
Gut-Instincts [1]	5eedev	0.8286	0.8480	<u>0.8382</u>
NLPatVCU [75]	ensemble1	0.8255	<u>0.8488</u>	0.8370
ICUE [38]	ensemble5-th10	<u>0.8369</u>	0.8294	0.8331
LYX-DMIP-FDU [44]	EnsembleBERT	0.8020	0.8513	0.8259
ata2425ds [NA]	transformer	0.7914	0.8432	0.8164
greenday [26]	llmner	0.7957	0.8278	0.8114
Graphswise-1 [16]	NERWise	0.8066	0.7955	0.8010
BASELINE [48]	NuNerZero-Finetuned	0.7639	0.8238	0.7927

Table 15. Performance metrics of each team’s top run beating the baseline for TB-RE. The best result is in bold, the second-best is underlined (micro-averaged).

Team ID	Run name	Precision	Recall	F1
Gut-Instincts [1]	6219eedev3re	0.6304	0.7532	0.6864
ONTUG [31]	ElectraCLEANR	0.7121	0.6104	<u>0.6573</u>
Graphswise-1 [16]	AtlopOnto	0.7418	0.5844	0.6538
ataupd2425-pam [62]	BiomedNLP-FULL_DATASET	0.5671	<u>0.7316</u>	0.6389
BIU-ONLP [32]	RobertaLarge	<u>0.7453</u>	0.5195	0.6122
BASELINE [48]	Atlop-Finetuned	0.7584	0.4892	0.5947

Table 16. Performance metrics of each team’s top run beating the baseline for TT-RE. The best result is in bold, the second-best is underlined (micro-averaged).

Team ID	Run name	Precision	Recall	F1
Gut-Instincts [1]	6229eedev3re	0.6280	<u>0.7572</u>	0.6866
ataupd2425-pam [62]	BiomedNLP-FULL_DATASET	0.5853	0.7202	<u>0.6458</u>
ONTUG [31]	ElectraCLEANR	0.7059	0.5926	0.6443
Graphswise-1 [16]	AtlopOnto	0.7326	0.5638	0.6372
ICUE [38]	biolinkbertl_pp	0.4974	0.7860	0.6093
BIU-ONLP [32]	RobertaLarge	<u>0.7362</u>	0.4938	0.5911
BASELINE [48]	Atlop-Finetuned	0.7533	0.4650	0.5751

BioNNE-L on nested named entity linking in Russian and English. (v) *Task ELCardioCC* on clinical coding in cardiology. (vi) *Task GutBrainIE* on gut-brain interplay information extraction.

The results for Task 13b suggest that the top participant systems achieved high scores, especially for yes/no answer generation, even in Phase A+, where no ground-truth relevant material was given. For list and factoid questions, system performance is less consistent, especially in Phase A+, indicating the presence of room for improvement. For these questions, the availability of ground-truth relevant material seems to allow the systems to provide answers of better quality.

Table 17. Performance metrics of each team’s top run beating the baseline for TM-RE. The best result is in bold, the second-best is underlined (micro-averaged).

Team ID	Run name	Precision	Recall	F1
Gut-Instincts [1]	6239eedev3re	0.4215	0.5147	0.4635
Graphswise-1 [16]	AtlopOnto	<u>0.4686</u>	0.3097	<u>0.3729</u>
ICUE [38]	biolinkbertl_pp	0.2858	<u>0.5054</u>	0.3651
LYX-DMIIP-FDU [44]	BioLinkBERT	0.3682	0.3257	0.3457
ONTUG [31]	ElectraCLEANR	0.3529	0.3231	0.3373
BASELINE [48]	Atlop-Finetuned	0.4986	0.2453	0.3288

This highlights the importance of Phase A, on the automated retrieval of relevant material, where the top performance is less consistent across batches, potentially affected by the domain of the expert posing the questions. A diverse set of retrieval and generation techniques was applied, including traditional methods, LLM-based frameworks, and integration of domain-specific knowledge. The results of task Synergy13, aligned with those of previous versions, suggest that state-of-the-art systems can be a useful tool for biomedical scientists in need of specialized information for developing problems, despite their limitations and room for improvement.

The new task MultiClinSum presented new challenging sub-tasks about text summarization of clinical case reports in Spanish, English, French and Portuguese. This task introduces the nuance of creating clinical Text Summarization systems specifically for the cardiology domain. In addition, it expands the range of the task beyond Spanish by introducing a sub-track that also involves English and Italian text. The results highlight the importance of having data specific to the language and specialty the systems are going to be applied in, even within domains that are already quite specific, like the clinical one.

The BioNNE-L task focused on the linking of biomedical entities for disorders, chemicals, and anatomical structures in Russian and English texts, addressing challenges such as nested entities and cross-language linking amid incomplete low-resource vocabularies. Despite the overall prevalence of LLMs in numerous domains and tasks, the top-performing systems for BioNNE-L utilized biomedical BERT-based retrieval and reranking architectures, highlighting the importance of task and domain-specific methods for information extraction.

The ELCardioCC task centered on extracting and classifying medical entities from Greek discharge letters, drawing participation from five teams who submitted a diverse set of systems across three subtasks. Most approaches leveraged transformer-based models—particularly BERT variants and multilingual LLMs—with strategies ranging from fine-tuned token classification to prompt-based extraction and embedding-based entity linking. The MLC-X subtask saw innovative uses of multilingual embeddings and LLM reasoning to predict and justify ICD-10 codes. Baseline models, built on multilingual BERT architectures, offered simple yet effective benchmarks for each subtask, emphasizing clarity and reproducibility.

The GutBrainIE task, centered on information extraction for the gut-brain axis, challenged participants with Named Entity Recognition and increasingly fine-grained Relation Extraction subtasks. Teams achieving the strongest performance employed supervised deep-learning strategies, combining pretrained biomedical language models with ensemble strategies. Only a few participants experimented with prompt-based or generative approaches; however, these generally obtained lower scores, confirming the need to develop specialized models to effectively extract complex entities and relations in a specific biomedical domain.

Overall, the participation in BioASQ 13 was significantly increased, both due to increased interest in the new versions of its already established tasks, as well as due to the introduction of four novel tasks. Several participating systems achieved competitive performance on the BioASQ tasks, and some of them managed to improve over the baselines or the state-of-the-art performance from previous years. Aligned with previous versions, BioASQ keeps pushing the research frontier in biomedical semantic indexing and question answering for thirteen years now, offering both well-established and new tasks. Initially, it extended beyond the English language and biomedical literature with the introduction of the task MESINESP [22] and continued consistently ever since. In this thirteenth edition, BioASQ was further extended with four new tasks, MultiClinSum [66], BioNNE-L [67], ElCardioCC [18], and GrutBrainIE [48]. As a result, BioASQ 13 offered tasks in six languages (English, Spanish, French, Portuguese, Russian, and Greek), three types of documents (biomedical articles, clinical case reports, and discharge letters), and two specialized domains within biomedicine (cardiology and gut-brain interaction).

The future directions for the BioASQ challenge involve further expanding the benchmark dataset for question answering through a community-driven approach, broadening the network of biomedical experts participating in the Synergy task, and enhancing the scope of resources used in the BioASQ tasks. This includes incorporating additional document types, multiple languages, and more specialized sub-domains within biomedicine.

6 Acknowledgments

The thirteenth edition of BioASQ is sponsored by Ovid, Atypion Systems Inc, and Elsevier. The MEDLINE/PubMed data resources considered in this work were accessed courtesy of the U.S. National Library of Medicine. BioASQ is grateful to the CMU team for providing the *exact answer* baselines for task 13b. This research was funded by the Ministerio de Ciencia e Innovación (MICINN) under project BARITONE (TED2021-129974B-C22). This work is also supported by the European Union’s Horizon Europe Co-ordination & Support Action under Grant Agreement No 101080430 (AI4HF), as well as Grant Agreement No 101057849 (DataTool4Heartproject). The work on the BioNNE-L task was supported by the Russian Science Foundation [grant number 23-11-00358]. ElCardioCC has been partially supported by project MIS 5154714 of the National Recovery and Resilience Plan Greece 2.0 funded by the European Union under

the NextGenerationEU Program. The work on the GutBrainIE task was supported by the HEREDITARY Project, as part of the European Union’s Horizon Europe research and innovation programme (GA 101137074).

References

1. Andersen, L.R., Gardshodn, M.I., Dolmer, M.H., Rodriguez, J.M., Dell’Aglio, D.: Trusting Gut Instincts: Transformer-Based Extraction of Structured Data from Gut-Brain Axis Publications. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
2. Angulo, J., Yeste, V.: AQAMS and AQAMS2: Multi Agent Systems for Biomedical Question Answering . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
3. Appleton, J.: The gut-brain axis: influence of microbiota on mood and mental health. *Integrative Medicine: A Clinician’s Journal* **17**(4), 28 (2018)
4. Ateia, S., Kruschwitz, U.: Can Language Models Critique Themselves? Investigating Self-Feedback for Retrieval Augmented Generation at BioASQ 2025 . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
5. Beltagy, I., Lo, K., Cohan, A.: SciBERT: Pretrained Language Model for Scientific Text. In: EMNLP (2019)
6. Bing-Chen, C., Han, J.C., Hung, H.C., Tsai, R.T.H.: NCU-IISR: Biomedical Question Answering via Gemini and GPT APIs in the BioASQ 13b Phase B Challenge . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
7. Bodenreider, O.: The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic acids research* **32**(suppl_1), D267–D270 (2004)
8. Bogdanov, S., Constantin, A., Bernard, T., Crabbé, B., Bernard, E.: NuNER: Entity Recognition Encoder Pre-training via LLM-Annotated Data (2024)
9. Borazio, F., Croce, D., Basili, R.: UniTor at BioASQ 2025: Modular Biomedical QA with Synthetic Snippets and Multiple Task Answer Generation . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
10. Burlova, A.: Navigating Partial UMLS Terminology: GAT Embeddings and Confidence Analysis for Multilingual Concept Linking. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
11. Carabotti, M., Scirocco, A., Maselli, M.A., Severi, C.: The gut-brain axis: interactions between enteric microbiota, central and enteric nervous systems. *Annals of gastroenterology: quarterly publication of the Hellenic Society of Gastroenterology* **28**(2), 203 (2015)
12. Clark, K., Luong, M.T., Le, Q.V., Manning, C.D.: ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. In: International Conference on Learning Representations (2020)
13. Conceição, S.I.R., Lopes, P.R.C., Couto, F.M.: lasigeBioTM at BioASQ25 Task GutBrainIE - Lean Large language models with syntactic features. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
14. Cryan, J.F., O’Riordan, K.J., Sandhu, K., Peterson, V., Dinan, T.G.: The gut microbiome in neurological disorders. *The Lancet Neurology* **19**(2), 179–194 (2020)

15. D. Lain, A., Lee, C., Doneva, S.E., Rodríguez-Cubillos, M.J., Castagnari, E., Simpson, T.I., , Posma, J.M.: Multilingual and Nested Biomedical Named Entity Normalisation via Candidate Retrieval and Lightweight Large Language Model Disambiguation. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
16. Datseris, A., Kuzmanov, M., Nikolova-Koleva, I., Taskov, D., Boytcheva, S.: Graph-wise @ CLEF-2025 GutBrainIE: Towards Automated Discovery of Gut-Brain Interactions: Deep Learning for NER and Relation Extraction from PubMed Abstracts. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
17. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp. 4171–4186. ACL, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>
18. Dimitriadis, D., Patsiou, V., Stoikopoulou, E., Toumpas, A., Kipouros, A., Papadopoulos, D., Bekiaridou, A., Barmpagiannos, K., Vasilopoulou, A., Barmpagiannos, A., Samaras, A., Giannakoulas, G., Tsoumakas, G.: Overview of ElCardioCC Task on Clinical Coding in Cardiology at BioASQ 2025. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
19. Dueñas Romero, S., Ureña-López, L.A., Martínez-Cámara, E.: SINAI at CLEF 2025: A Multi-Stage RAG Pipeline for Biomedical Semantic Question Answering . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
20. Feng, F., Yang, Y., Cer, D., Arivazhagan, N., Wang, W.: Language-agnostic BERT sentence embedding. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 878–891. ACL, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.acl-long.62>
21. Galat, D., Molla-Aliod, D.: LLM Ensemble for RAG: Role of context length in zero-shot Question Answering for BioASQ Challenge . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
22. Gasco, L., Nentidis, A., Krithara, A., Estrada-Zavala, D., Toshiyuki Murasaki, R., Primo-Peña, E., Bojo-Canales, C., Paliouras, G., Krallinger, M.: Overview of BioASQ 2021-MESINESP track. Evaluation of advance hierarchical classification techniques for scientific literature, patents and clinical trials. In: Proceedings of the 9th BioASQ Workshop (2021)
23. Ghaisas, S., Maher, J., Kanthasamy, A.: Gut microbiome in health and disease: Linking the microbiome–gut–brain axis and environmental factors in the pathogenesis of systemic and neurodegenerative diseases. *Pharmacology & therapeutics* **158**, 52–62 (2016)
24. Grazhdanski, G.: Group relative policy optimization for spanish clinical case report summarization. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
25. Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., Poon, H.: Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)* **3**(1), 1–23 (2021)
26. Gupta, H.P., Banerjee, R.: LLMs for Biomedical NER. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)

27. Han, J., Liu, Y.: GutUZH at CLEF2025 BioASQ Task 6: a method of SOTA performance with the best results at GutBrainIE NER subtask 1. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
28. Huang, B.: Clinical entity recognition and linking in greek discharge letters using multilingual-llm-based multi-stage system. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
29. Huguet Cabot, P.L., Navigli, R.: REBEL: Relation extraction by end-to-end language generation. In: Findings of the Association for Computational Linguistics: EMNLP 2021. pp. 2370–2381. ACL, Punta Cana, Dominican Republic (Nov 2021), <https://aclanthology.org/2021.findings-emnlp.204>
30. Jonker, R.A.A., Almeida, T., Almeida, J., Matos, S.: BIT.UA at BioASQ 13B: Revisiting Evaluation, DPRF-Enhanced Retrieval and Fine-Tuned LLMs. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
31. Kantz, B., Waldert, P., Lengauer, S., Schreck, T.: Constrained Linked Entity Annotation using RAG (CLEANR). In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
32. Keinan, R., Cohen, A.D.N., Tsarfaty, R.: From Named Entities to Relations: End-to-End Biomedical Information Extraction. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
33. Kim, H., Lee, H., Cho, Y., Park, J., Park, J., Park, S., Chok, Y.T., Baek, S., Lee, D., Kang, J.: Prompting Matters: Snippet-Aware Strategies for Biomedical QA with LLMs in BioASQ 13b . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
34. Krithara, A., Nentidis, A., Bougiatiotis, K., Paliouras, G.: BioASQ-QA: A manually curated corpus for Biomedical Question Answering. *Scientific Data* **10**(1), 170 (2023)
35. Krithara, A., Nentidis, A., Paliouras, G., Kakadiaris, I.: Results of the 4th edition of BioASQ Challenge. In: Proceedings of the Fourth BioASQ workshop (2016), <https://www.aclweb.org/anthology/W16-3101.pdf>
36. Krithara, A., Nentidis, A., Paliouras, G., Krallinger, M., Miranda, A.: BioASQ at CLEF2021: large-scale biomedical semantic indexing and question answering. In: Advances in Information Retrieval: ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II 43. pp. 624–630. Springer (2021)
37. Krithara, A., Nentidis, A., Vantorou, E., Katsimpras, G., Almirantis, Y., Arnal, M., Bunevicius, A., Farre-Maduell, E., Kassiss, M., Konstantakos, V., Matis-Mitchell, S., Polychronopoulos, D., Rodriguez-Pascual, J., Samaras, E.G., Samiotaki, M., Sanoudou, D., Vozi, A., Paliouras, G.: BioASQ Synergy: a dialogue between question-answering systems and biomedical experts for promoting COVID-19 research. *Journal of the American Medical Informatics Association* p. ocae232 (08 2024). <https://doi.org/10.1093/jamia/ocae232>
38. Lee, C., Doneva, S., Rodriguez-Cubillos, M., Castagnari, E., Lain, A., Posma, J., Simpson, T.I.: Understanding Gut-Brain Interplay in Scientific Literature: A Hybrid Approach from Classification to Generative LLM Reasoning. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
39. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* **36**(4), 1234–1240 (09 2019). <https://doi.org/10.1093/bioinformatics/btz682>
40. Li, C., Zheng, X., Liu, S.: BIBERT on Biomedical Nested Named Entity Linking at BioASQ 2025. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)

41. Lin, C.Y.: ROUGE: A package for automatic evaluation of summaries. In: Proceedings of the ACL workshop ‘Text Summarization Branches Out’. pp. 74–81. Barcelona, Spain (2004)
42. Liu, F., Shareghi, E., Meng, Z., Basaldella, M., Collier, N.: Self-alignment pre-training for biomedical entity representations. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 4228–4238. ACL, Online (Jun 2021). <https://doi.org/10.18653/v1/2021.naacl-main.334>
43. Liu, F., Vulić, I., Korhonen, A., Collier, N.: Learning domain-specialised representations for cross-lingual biomedical entity linking. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers). pp. 565–574. ACL, Online (Aug 2021). <https://doi.org/10.18653/v1/2021.acl-short.72>, <https://aclanthology.org/2021.acl-short.72/>
44. Liu, Y.: LYX_DMIP_FDU at BioASQ 2025: Utilizing BERT embeddings for biomedical text mining. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
45. Loukachevitch, N., Manandhar, S., Baral, E., Rozhkov, I., Braslavski, P., Ivanov, V., Batura, T., Tutubalina, E.: NEREL-BIO: A Dataset of Biomedical Abstracts Annotated with Nested Named Entities. *Bioinformatics* (04 2023). <https://doi.org/10.1093/bioinformatics/btad161>, btad161
46. Loukachevitch, N., Sakhovskiy, A., Tutubalina, E.: Biomedical concept normalization over nested entities with partial UMLS terminology in Russian. In: Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 2383–2389. ELRA and ICCL, Torino, Italia (May 2024), <https://aclanthology.org/2024.lrec-main.213/>
47. Malakasiotis, P., Pavlopoulos, I., Androutsopoulos, I., Nentidis, A.: Evaluation measures for task b. Tech. rep., Tech. rep. BioASQ (2022), http://participants-area.bioasq.org/Tasks/b/eval_meas.2022
48. Martinelli, M., Silvello, G., Bonato, V., Di Nunzio, G.M., Ferro, N., Irrera, O., Marchesin, S., Menotti, L., Vezzani, F.: Overview of GutBrainIE@CLEF 2025: Gut-Brain Interplay Information Extraction. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
49. Mehta, R.: Enhancing Biomedical Named Entity Recognition using GLiNER-BioMed with Targeted Dictionary-Based Post-processing for BioASQ 2025 task 6. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
50. Nentidis, A., Katsimpras, G., Krithara, A., Krallinger, M., Ortega, M.R., Loukachevitch, N., Sakhovskiy, A., Tutubalina, E., Tsoumakas, G., Giannakoulas, G., Bekiaridou, A., Samaras, A., Di Nunzio, G.M., Ferro, N., Marchesin, S., Menotti, L., Silvello, G., Paliouras, G.: BioASQ at CLEF2025: The Thirteenth Edition of the Large-Scale Biomedical Semantic Indexing and Question Answering Challenge. In: *Advances in Information Retrieval*. pp. 407–415. Springer Nature Switzerland, Cham (2025)
51. Nentidis, A., Katsimpras, G., Krithara, A., Lima López, S., Farré-Maduell, E., Gasco, L., Krallinger, M., Paliouras, G.: Overview of BioASQ 2023: The Eleventh BioASQ Challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. In: Arampatzis, A., Kanoulas, E., Tsikrika, T., Vrochidis, S., Giachanou,

- A., Li, D., Aliannejadi, M., Vlachos, M., Faggioli, G., Ferro, N. (eds.) *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. pp. 227–250. Springer Nature Switzerland, Cham (2023)
52. Nentidis, A., Katsimpras, G., Krithara, A., Lima-López, S., Farré-Maduell, E., Krallinger, M., Loukachevitch, N., Davydova, V., Tutubalina, E., Paliouras, G.: Overview of BioASQ 2024: The twelfth BioASQ challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Fifteenth International Conference of the CLEF Association (CLEF 2024)* (2024)
 53. Nentidis, A., Katsimpras, G., Krithara, A., Paliouras, G.: Overview of BioASQ Tasks 11b and Synergy11 in CLEF2023. In: *CEUR Workshop Proceedings* (2023)
 54. Nentidis, A., Katsimpras, G., Krithara, A., Paliouras, G.: Overview of BioASQ Tasks 12b and Synergy12 in CLEF2024. In: Faggioli, G., Ferro, N., Galuščáková, P., García Seco de Herrera, A. (eds.) *CLEF Working Notes* (2024)
 55. Nentidis, A., Katsimpras, G., Krithara, A., Paliouras, G.: Overview of BioASQ Tasks 13b and Synergy13 in CLEF2025. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 56. Nentidis, A., Katsimpras, G., Vondrou, E., Krithara, A., Gasco, L., Krallinger, M., Paliouras, G.: Overview of BioASQ 2021: The Ninth BioASQ Challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. In: *International Conference of the Cross-Language Evaluation Forum for European Languages*. pp. 239–263. Springer (2021)
 57. Nentidis, A., Katsimpras, G., Vondrou, E., Krithara, A., Miranda-Escalada, A., Gasco, L., Krallinger, M., Paliouras, G.: Overview of BioASQ 2022: The Tenth BioASQ Challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. Springer (2022). https://doi.org/10.1007/978-3-031-13643-6_22
 58. Nentidis, A., Katsimpras, G., Vondrou, E., Krithara, A., Paliouras, G.: Overview of BioASQ Tasks 9a, 9b and Synergy in CLEF2021. In: *Proceedings of the 9th BioASQ Workshop A challenge on large-scale biomedical semantic indexing and question answering*. CEUR Workshop Proceedings (2021), <http://ceur-ws.org/Vol-2936/paper-10.pdf>
 59. Nentidis, A., Katsimpras, G., Vondrou, E., Krithara, A., Paliouras, G.: Overview of BioASQ Tasks 10a, 10b and Synergy10 in CLEF2022. In: *CEUR Workshop Proceedings*. vol. 3180, pp. 171–178 (2022)
 60. Nentidis, A., Krithara, A., Paliouras, G., Krallinger, M., Sanchez, L.G., Lima, S., Farre, E., Loukachevitch, N., Davydova, V., Tutubalina, E.: BioASQ at CLEF2024: The Twelfth Edition of the Large-Scale Biomedical Semantic Indexing and Question Answering Challenge. In: *ECIR2024*. pp. 490–497. Springer (2024)
 61. van den Oord, A., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. *ArXiv* **abs/1807.03748** (2018), <https://api.semanticscholar.org/CorpusID:49670925>
 62. Pamio, L., Di Nunzio, G.M.: BioASQ task GutBrainIE 2025 Task 6: Comparing CRF vs BERT Models for Named Entity Recognition and Relation Extraction. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 63. Panou, D., Dimopoulos, A., Koubarakis, M., Reczko, M.: Harnessing Collective Intelligence of LLMs for Robust Biomedical QA: A Multi-Model Approach. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)

64. Peña Gnecco, D., Serrano, J., Puertas, E., Martínez-Santos, J.C.: Hybrid Re-ranking for Biomedical Entity Linking using SapBERT Embeddings: A High-Performance System for BioNNE-L 2025-1. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
65. Piron, S., Di Nunzio, G.M.: Named Entity Recognition with GLiNER and Relation Extraction with LLMs. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
66. Rodríguez-Ortega, M., Rodríguez-Lopez, E., Lima-López, S., Escolano, C., Melero, M., Pratesi, L., Vigil-Gimenez, L., Fernandez, L., Farré-Maduell, E., Krallinger, M.: Overview of MultiClinSum task at BioASQ 2025: evaluation of clinical case summarization strategies for multiple languages: data, evaluation, resources and results. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
67. Sakhovskiy, A., Loukachevitch, N., Tutubalina, E.: Overview of the BioASQ BioNNE-L Task on Biomedical Nested Entity Linking in CLEF 2025. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
68. Sakhovskiy, A., Semenova, N., Kadurin, A., Tutubalina, E.: Graph-enriched biomedical entity representation transformer. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. pp. 109–120. Springer Nature Switzerland, Cham (2023)
69. Sakhovskiy, A., Semenova, N., Kadurin, A., Tutubalina, E.: Biomedical entity representation with graph-augmented multi-objective transformer. In: *Findings of the Association for Computational Linguistics: NAACL 2024*. pp. 4626–4643. ACL, Mexico City, Mexico (Jun 2024). <https://doi.org/10.18653/v1/2024.findings-naacl.288>
70. Salton, G., Buckley, C.: Improving retrieval performance by relevance feedback. *Journal of the American Society for Information Science* **41**(4), 288–297 (jun 1990). [https://doi.org/10.1002/\(SICI\)1097-4571\(199006\)41:4<288::AID-ASIS4.3.0.CO;2-H](https://doi.org/10.1002/(SICI)1097-4571(199006)41:4<288::AID-ASIS4.3.0.CO;2-H)
71. Schneider, E.T.R., Schneider, F.H., Paraiso, E.C., Britto Jr, A.S., Cruz, R.M.O.: MedGemma-Sum-Pt: A Lightweight Model for Portuguese Clinical Summarization. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
72. Stachura, D., Konieczna, J., Nowak, A.: Are Smaller Open-Weight LLMs Closing the Gap to Proprietary Models for Biomedical Question Answering? . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
73. Sung, M., Jeon, H., Lee, J., Kang, J.: Biomedical entity representations with synonym marginalization. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 3641–3650. ACL, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.335>
74. Tang, J., Yang, H., Xiong, K., Li, H., Quaresma, P., Yu, H., Zhang, W., Song, M., Jiang, Y.: Applying DeepSeek to BioASQ Task 13B: Using Supervised Fine-Tuning and Few-Shot Learning . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
75. Taylor, S., Dil, C., Shah, A., Jannat, Oldham, C., Upadhyay, A., Varughese, J., Yazbeck, N., McInnes, B.T.: NLP@VCU at BioASQ2025: Information Extraction on the GutBrainIE dataset. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) CLEF 2025 Working Notes (2025)
76. Tedeschi, S., Maiorca, V., Campolungo, N., Cecconi, F., Navigli, R.: WikiNEu-Ral: Combined neural and knowledge-based silver data creation for multilingual NER. In: *Findings of the Association for Computational Linguistics: EMNLP*

2021. pp. 2521–2533. ACL, Punta Cana, Dominican Republic (Nov 2021). <https://doi.org/10.18653/v1/2021.findings-emnlp.215>
77. Tsatsaronis, G., Balikas, G., Malakasiotis, P., Partalas, I., Zschunke, M., Alvers, M.R., Weissenborn, D., Krithara, A., Petridis, S., Polychronopoulos, D., Almirantis, Y., Pavlopoulos, J., Baskiotis, N., Gallinari, P., Artieres, T., Ngonga, A., Heino, N., Gaussier, E., Barrio-Alvers, L., Schroeder, M., Androutsopoulos, I., Paliouras, G.: An overview of the BIOASQ large-scale biomedical semantic indexing and question answering competition. *BMC Bioinformatics* **16**, 138 (2015)
 78. Vachharajani, P.: Multilingual embedding and prompt-driven approaches for named entity recognition, entity linking, and clinical code prediction in greek discharge summaries. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 79. Vachharajani, P.: pmathematician at MultiClinSUM 2025: A Novel Automated Prompt Optimization Framework for Multilingual Clinical Summarization. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 80. Velichkov, B., Datseris, A., Vassileva, S., Boytcheva, S.: Enigma @ ElCardioCC: Bridging NER and ICD-10 Entity Linking - A Hybrid Method for Greek Clinical Narratives. In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 81. Verma, S., Jiang, F., Xue, X.: Beyond Retrieval: Ensembling Cross-Encoders and GPT Rerankers with LLMs for Biomedical QA . In: Faggioli, G., Ferro, N., Rosso, P., Spina, D. (eds.) *CLEF 2025 Working Notes* (2025)
 82. Yang, Z., Zhou, Y., Nyberg, E.: Learning to answer biomedical questions: OAQA at BioASQ 4B. In: Kakadiaris, I.A., Paliouras, G., Krithara, A. (eds.) *Proceedings of the Fourth BioASQ workshop*. pp. 23–37. ACL, Berlin, Germany (Aug 2016). <https://doi.org/10.18653/v1/W16-3104>
 83. Yasunaga, M., Leskovec, J., Liang, P.: LinkBERT: Pretraining Language Models with Document Links. In: *Association for Computational Linguistics (ACL)* (2022)
 84. Zaratiana, U., Tomeh, N., Holat, P., Charnois, T.: GLiNER: Generalist model for named entity recognition using bidirectional transformer. In: Duh, K., Gomez, H., Bethard, S. (eds.) *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. pp. 5364–5376. ACL, Mexico City, Mexico (Jun 2024). <https://doi.org/10.18653/v1/2024.naacl-long.300>
 85. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: BERTScore: Evaluating Text Generation with BERT. In: *International Conference on Learning Representations (ICLR)* (2020), <https://arxiv.org/abs/1904.09675>
 86. Zhou, W., Huang, K., Ma, T., Huang, J.: Document-level relation extraction with adaptive thresholding and localized context pooling. In: *Proceedings of the AAAI Conference on Artificial Intelligence* (2021)
 87. Zhuang, L., Wayne, L., Ya, S., Jun, Z.: A robustly optimized BERT pre-training approach with post-training. In: Li, S., Sun, M., Liu, Y., Wu, H., Liu, K., Che, W., He, S., Rao, G. (eds.) *Proceedings of the 20th Chinese National Conference on Computational Linguistics*. pp. 1218–1227. Chinese Information Processing Society of China, Huhhot, China (Aug 2021), <https://aclanthology.org/2021.ccl-1.108/>