

**Algorithm 5:** WCKMeans**Input:** Corpus  $\mathcal{T}$ , clusters  $K$ , MLS  $\mathcal{M}$ , CLS  $\mathcal{C}$ , max iterations  $T$ **Output:** Labels  $\Gamma$ 


---

```

1 Normalize  $\mathcal{T}$  to get embeddings  $\mathbf{X}$ ;
2 Construct constraint matrix  $\mathbf{R}$  with  $\mathcal{M}$ ,  $\mathcal{C}$ ,  $|\mathcal{T}|$ , and  $\mathbf{X}$ ;
3 Initialize centers  $\mathcal{U} \leftarrow \{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ ;
4 while not converged &  $t < T$  do
5   foreach  $x_i \in \mathbf{X}$  do
6      $\Gamma_i \leftarrow \arg \min_{1 \leq k \leq K} \|x_i - \mathbf{u}_k\|^2 + \sum_{\substack{(i,j) \in \mathcal{M} \\ \Gamma_j \neq k}} R_{i,j} + \sum_{\substack{(i,j) \in \mathcal{C} \\ \Gamma_j = k}} |R_{i,j}|$ 
7   Update centers  $\mathbf{u}_k \leftarrow \text{Mean}(\{x_i \mid \Gamma_i = k\})$  for each  $k$ ;
8 return  $\Gamma$ 

```

---

**A Proof of Lemma 5.2**

PROOF. Let  $(t_i, t_j) = \arg \max_{(t_a, t_b) \in S, a < b} a$ . In the extreme case,  $j = i + 1$  and only edges  $(t_a, t_{a+1}), (t_a, t_{a+2}), \dots, (t_a, t_j)$  are added into  $S$  for each  $a \in [1, i - 1]$  as they are surely larger than  $(t_i, t_j)$ . Put in another way, for each  $t_a$  with  $a \in [1, i]$ , we can select  $t_b$  with  $a + 1 \leq b \leq j$ , i.e., which corresponds to  $j - a = i + 1 - a$  text instances. Thus, the total number of edges satisfy:

$$1 + 2 + \dots + i = \frac{(i+1)i}{2} = N.$$

Then, we can derive

$$i^2 + i - 2N = 0 \Rightarrow i = \frac{\sqrt{8N+1}-1}{2},$$

which finishes the proof.  $\square$

**B Details of Datasets and Baselines**

**Datasets.** The *BBC News* [27] dataset contains 2,225 headlines collected from the BBC News website and categorized into 5 high-level topics. It serves as a benchmark for multi-class classification tasks in the news field. The *Tweet* [71] dataset consists of 2,472 tweets annotated for relevance to 89 queries from the TREC Microblog track, and is widely used for short-text retrieval and ranking. The *Bank77* [74] dataset includes 3,080 customer service utterances mapped to 77 intent categories, enabling fine-grained intent classification in the financial category. The *Reddit* [36] dataset, with 3,217 posts from online communities, is used primarily for unsupervised topic discovery. The *CLINC* [78] dataset comprises 4,500 user requests labeled across 10 distinct domains for domain-level text clustering. Finally, the *Massive* [20] dataset contains 11,514 utterances spanning 18 scenario-based categories, and is often used in multilingual and multi-scenario classification settings.

**Baselines.** We include two sets of baselines in the experiments: *embedding-based clustering methods* and *LLM-assisted clustering methods*. In embedding-based clustering methods, we include seven representative baselines that follow a two-stage paradigm: generating fixed-dimensional text embeddings followed by conventional clustering algorithms, such as K-Means++ or *spectral clustering* [63]. The selected embedding models encompass both traditional and neural models.

- **TF-IDF:** A classic sparse representation based on term frequency-inverse document frequency, widely used in earlier text clustering studies.
- **E5** [65]: A recent retrieval-oriented embedding model pretrained with a multi-task objective, shown to perform well in clustering tasks.
- **DistilBERT** and **Sentence-BERT** [46]: Popular lightweight transformers optimized for sentence-level embeddings.
- **Instructor-Large** [56]: A powerful instruction-tuned model designed to align embeddings with various downstream tasks; we use it as the default embedder in most baselines due to its strong performance.
- **OpenAI-GPT** [6] and **LLaMA-2 (7B)** [60]: Large language models (LLMs) whose embeddings are obtained by averaging final-layer token representations. As LLMs are often optimized for generation, we treat them as black-box embedding extractors here.

These embedding models are coupled with K-Means++ or *spectral clustering*, depending on the context.

In LLM-assisted clustering methods, we additionally consider three recent strong methods that leverage LLMs beyond embedding extraction.

- **SCCL** [73]: A self-supervised contrastive clustering method that improves cluster discrimination by minimizing intra-cluster distance and maximizing inter-cluster margins in embedding space.
- **ClusterLLM** [77]: Utilizes LLMs to provide clustering feedback and refines smaller embedding models through LLM guidance. We also include a variant that excludes fine-tuning to assess the impact of feedback alone.
- **PO-PCKMeans** [62]: A semi-supervised approach that constructs pairwise constraints from an LLM based on a small set of labeled examples and integrates them into a constrained clustering framework.

**C Additional Algorithmic Details****C.1 WCKMeans**

Algorithm 5 integrates well-designed weighted constraints into standard K-Means by creating a constraint matrix  $\mathbf{R}$ , where  $R_{i,j} > 0$  for must-links and  $R_{i,j} < 0$  for cannot-links. In each iteration, each point  $x_i$  is assigned to the cluster  $k$  that minimizes its squared distance to centroid  $\mathbf{u}_k$  plus penalties: add  $R_{i,j}$  if a must-link partner  $j$  is in a different cluster, and add  $|R_{i,j}|$  if a cannot-link partner  $j$  is in the same cluster. After assignments, centroids are updated as the mean of their assigned points. This repeats until convergence or  $T$  iterations.

**C.2 WCSC**

Algorithm 6 starts by building a graph that captures how similar each pair of points is in the normalized embedding space. Then selecting top- $K$  eigenvectors finds a low-dimensional embedding that balances the natural data structure with must-links and cannot-links. The  $\alpha$  [66] serves as a balance factor between the original information captured by the graph and the constraints. Finally, K-Means is applied to produce the cluster labels.

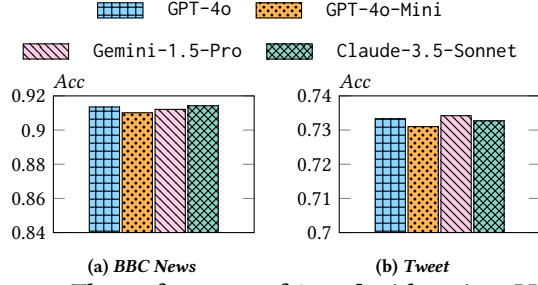


Figure 6: The performance of Cequel with various LLMs

**Algorithm 6: W CSC**

**Input:** Corpus  $\mathcal{T}$ , clusters  $K$ , MLS  $\mathcal{M}$ , CLS  $C$   
**Output:** Labels  $\Gamma$

- 1 Normalize  $\mathcal{T}$  to get embeddings  $X$ ;
- 2 Build affinity matrix  $A$  with  $X$ ;
- 3 Construct constraint matrix  $R$  with  $\mathcal{M}$ ,  $C$ ,  $|\mathcal{T}|$ , and  $X$ ;
- 4  $L \leftarrow D^{-1/2}(D - A)D^{-1/2}$ ,  $R_\alpha \leftarrow D^{-1/2}RD^{-1/2} - \alpha I$ ;
- 5 Compute top- $K$  generalized eigenvectors  $V$  of  $L$  and  $R_\alpha$ ;
- 6  $\Gamma \leftarrow K\text{-Means}(D^{-1/2}V)$ ;
- 7 **return**  $\Gamma$

**Algorithm 7: Edge Query**

**Input:** Corpus  $\mathcal{T}$ , the number of edges  $N$ , Oracle  $O$   
**Output:** MLS  $\mathcal{M}$ , CLS  $C$

- 1  $S \leftarrow \text{GreedyEdgeSelection}(\mathcal{T}, N)$ ;
- 2 **foreach**  $(t_i, t_j) \in S$  **do**
- 3     **switch**  $O(t_i, t_j)$  **do**
- 4         **case**  $ML$  **do**
- 5              $\mathcal{M} \leftarrow \mathcal{M} \cup \{(t_i, t_j)\}$
- 6         **case**  $CL$  **do**
- 7              $C \leftarrow C \cup \{(t_i, t_j)\}$
- 8 **return**  $\mathcal{M}, C$

**Algorithm 8: Triangle Query**

**Input:** Corpus  $\mathcal{T}$ , the number of triangles  $N_\Delta$ , Oracle  $O$   
**Output:** MLS  $\mathcal{M}$ , CLS  $C$

- 1  $S_\Delta \leftarrow \text{GreedyTriangleSelection}(\mathcal{T}, N_\Delta)$ ;
- 2 **foreach**  $(t_i, t_j, t_k) \in S_\Delta$  **do**
- 3     **switch**  $O(t_i, t_j, t_k)$  **do**
- 4         **case**  $all\text{-}same$  **do**
- 5              $\forall (p, q) \in \{(i, j), (i, k), (j, k)\} : \mathcal{M} \leftarrow \mathcal{M} \cup \{(t_p, t_q)\}$
- 6         **case**  $ij\text{-}same$  **do**
- 7              $\mathcal{M} \leftarrow \mathcal{M} \cup \{(t_i, t_j)\}$ ,  $C \leftarrow C \cup \{(t_i, t_k), (t_j, t_k)\}$
- 8         **case**  $ik\text{-}same$  **do**
- 9              $\mathcal{M} \leftarrow \mathcal{M} \cup \{(t_i, t_k)\}$ ,  $C \leftarrow C \cup \{(t_i, t_j), (t_j, t_k)\}$
- 10         **case**  $jk\text{-}same$  **do**
- 11              $\mathcal{M} \leftarrow \mathcal{M} \cup \{(t_j, t_k)\}$ ,  $C \leftarrow C \cup \{(t_i, t_j), (t_i, t_k)\}$
- 12         **case**  $all\text{-}diff$  **do**
- 13              $\forall (p, q) \in \{(i, j), (i, k), (j, k)\} : C \leftarrow C \cup \{(t_p, t_q)\}$
- 14 **return**  $\mathcal{M}, C$

**C.3 Edge Query**

Algorithm 7 uses Algorithm 1 (GreedyEdgeSelection) to select the top  $N$  instance pairs most informative for querying. Each chosen pair  $(t_i, t_j)$  is queried to the oracle  $O$ , and the corresponding constraints are added based on the responses.

**C.4 Triangle Query**

Algorithm 8 utilizes Algorithm 4 (GreedyTriangleSelection) to select  $N_\Delta$  triangles  $(t_i, t_j, t_k)$ . For each triangle, the oracle  $O$  returns five outcomes—"all-same", "ij-same", "ik-same", "jk-same", or "all-diff", and the corresponding must-link or cannot-link constraints are added to  $\mathcal{M}$  and  $C$ .

**D Additional Experiments**

**Various Text Encoders.** Table 7 reports the clustering performance of Cequel and strong baselines using two representative encoders: Sentence-BERT and Instructor-Large. Results demonstrate the consistent superiority of Cequel across datasets and embeddings. For instance, with Instructor-Large on *Tweet*, Cequel reaches 73.10% ACC and 89.08% NMI, exceeding ClusterLLM by 7.29% and 0.47%. With Sentence-BERT on *BBC News*, Cequel achieves 84.45% ACC and 61.65% NMI, outperforming PO-PCKMeans by 1.75% and 4.22%, respectively. These gains are consistent across diverse corpora, from structured news to informal tweets and task-oriented texts like *Bank77* and *Reddit*. The results highlight the adaptability of Cequel and its ability to take advantage of powerful pre-trained encoders and informative constraints without task-specific tuning.

**Queried LLMs.** Fig. 6 illustrates the clustering results of Cequel when querying four popular LLMs including GPT-4o, GPT-4o-Mini, Gemini-1.5-Pro, and Claude-3.5-Sonnet. Key observations include: (1) Minimal differences in the adoption of various LLMs. (2) LLMs influence answer accuracy, particularly regarding constraints, without affecting query triangle selection; for example, GPT-4o enhances constraint accuracy compared to GPT-4o-Mini, resulting in increased overall accuracy. (3) The advantages of more expensive LLMs do not warrant their costs, underscoring the necessity for more cost-effective models capable of managing complex tasks like triangle queries.

**E Additional Prompts**

We provide additional example prompts in Table 8 for edge and triangle queries across multiple datasets, including *BBC News*, *Tweet*, *Bank77*, *Reddit*, and *Massive*. Edge queries ask whether two documents belong to the same category (Yes/No), while triangle queries assess the categorical relationship among three documents via a five-way multiple-choice format. These prompts are designed to support constraint construction using LLMs for clustering tasks.

**Table 7: Clustering performance with various text encoders**

Encoder	Method	BBC News		Tweet		Bank77		Reddit	
		ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
Sentence-BERT	SCCL	82.00	<u>58.33</u>	53.60	72.80	27.80	53.21	36.60	45.71
	ClusterLLM	80.08	<u>54.13</u>	57.47	83.74	<u>58.61</u>	77.06	36.09	44.74
	PO-PCKMeans	<u>82.70</u>	57.43	<u>61.39</u>	<u>86.12</u>	57.73	<u>78.09</u>	<u>38.48</u>	<u>48.56</u>
	Ceque1	<b>84.45</b>	<b>61.65</b>	<b>78.48</b>	<b>89.30</b>	<b>67.05</b>	<b>80.46</b>	<b>40.94</b>	<b>49.32</b>
Instructor-Large	SCCL	83.60	62.00	36.50	68.50	36.00	61.30	28.70	36.00
	ClusterLLM	88.40	69.73	<u>65.81</u>	<u>88.61</u>	<u>68.27</u>	<u>82.45</u>	54.80	62.47
	PO-PCKMeans	<u>89.53</u>	<u>72.16</u>	65.58	88.20	67.79	81.92	<u>55.10</u>	<u>62.51</u>
	Ceque1	<b>91.01</b>	<b>75.22</b>	<b>73.10</b>	<b>89.08</b>	<b>69.03</b>	<b>83.20</b>	<b>57.78</b>	<b>64.25</b>

**Table 8: Example prompts for edge and triangle queries on different datasets.**

Dataset	Type	Prompt
BBC News	Edge Query	Cluster BBC News docs by whether they belong to the same news category. For each pair, respond with Yes or No without explanation. - News #1: Ad sales boost Time Warner profit - News #2: Air passengers win new EU rights Given this context, do News #1 and News #2 likely correspond to the same news category?
	Triangle Query	Cluster BBC News docs by whether they belong to the same news category. For each triangle, respond with a, b, c, d, or e without explanation. - News #1: REM announce new Glasgow concert - News #2: Moreno debut makes Oscar mark - News #3: Last Star Wars 'not for children' Given this context, do News #1, News #2, and News #3 likely correspond to the same news category? (a) All are same category. (b) Only #1 and #2 are same category. (c) Only #1 and #3 are same category. (d) Only #2 and #3 are same category. (e) None.
Tweet	Edge Query	Cluster Tweet docs by whether they belong to the same tweet category. For each pair, respond with Yes or No without explanation. - Tweet #1: super bowl commercial - Tweet #2: kung pao chicken recipe Given this context, do Tweet #1 and Tweet #2 likely correspond to the same tweet category?
	Triangle Query	Cluster Tweet docs by whether they belong to the same tweet category. For each triangle, respond with a, b, c, d, or e without explanation. - Tweet #1: weight loss diet plan acai juice equal nutritional claim acai - Tweet #2: watch christina aguilera screw national anthem super bowl post - Tweet #3: yell cast jr hartley ad digital era author search fly fishing book dj hunt book read Given this context, do Tweet #1, Tweet #2, and Tweet #3 likely correspond to the same tweet category? (a) All are same category. (b) Only #1 and #2 are same category. (c) Only #1 and #3 are same category. (d) Only #2 and #3 are same category. (e) None.
Bank77	Edge Query	Cluster Bank77 docs by whether they belong to the same intent category. For each pair, respond with Yes or No without explanation. - Intent #1: I think something went wrong with my card delivery as I haven't received it yet. - Intent #2: How do I link to my credit card with you? Given this context, do Intent #1 and Intent #2 likely correspond to the same intent category?
	Triangle Query	Cluster Bank77 docs by whether they belong to the same intent category. For each triangle, respond with a, b, c, d, or e without explanation. - Intent #1: What is the current exchange rate for me? - Intent #2: I believe my card payment exchange rate is incorrect. - Intent #3: I think my child used my card while I wasn't home. Given this context, do Intent #1, Intent #2, and Intent #3 likely correspond to the same intent category? (a) All are same category. (b) Only #1 and #2 are same category. (c) Only #1 and #3 are same category. (d) Only #2 and #3 are same category. (e) None.
Reddit	Edge Query	Cluster Reddit docs by whether they belong to the same topic category. For each pair, respond with Yes or No without explanation. - Topic #1: CTA to make all rail stations accessible within 20 years - Topic #2: A Krazy Mug Oyee Balle Balle!!! Kettle with set of two glasses Given this context, do Topic #1 and Topic #2 likely correspond to the same topic category?
	Triangle Query	Cluster Reddit docs by whether they belong to the same topic category. For each triangle, respond with a, b, c, d, or e without explanation. - Topic #1: Batman and Batwoman Fan art - Topic #2: Books on American Union and Labour Movements? - Topic #3: When Taliban offer you gold: Afghan youth in crisis? Given this context, do Topic #1, Topic #2, and Topic #3 likely correspond to the same topic category? (a) All are same category. (b) Only #1 and #2 are same category. (c) Only #1 and #3 are same category. (d) Only #2 and #3 are same category. (e) None.
Massive	Edge Query	Cluster Massive docs by whether they belong to the same scenario category. For each pair, respond with Yes or No without explanation. - Scenario #1: here is something from today - Scenario #2: remind me of the meeting on tuesday Given this context, do Scenario #1 and Scenario #2 likely correspond to the same scenario category?
	Triangle Query	Cluster Massive docs by whether they belong to the same scenario category. For each triangle, respond with a, b, c, d, or e without explanation. - Scenario #1: alexa play my country playlist - Scenario #2: write email to my company mate to submit the task tomorrow - Scenario #3: remind me about my business meeting at three and forty five p. m. Given this context, do Scenario #1, Scenario #2, and Scenario #3 likely correspond to the same scenario category? (a) All are same category. (b) Only #1 and #2 are same category. (c) Only #1 and #3 are same category. (d) Only #2 and #3 are same category. (e) None.