

Pemetaan Topik Pembicaraan Pada Komentar *Live Youtube* Menggunakan *K-Means Clustering* sebagai Identifikasi awal Kejahatan Verbal *Cyberbullying*

Alfirna Rizqi Lahitani^{1*}, Adlia Nur Zhafarina², Nanda Saputri Winda Oktavia³, Nita Jariyah⁴

^{1,3,4} Teknologi Informasi, Fakultas Teknik dan Teknologi Informasi Universitas Jenderal Achmad Yani Yogyakarta

² Hukum, Fakultas Ekonomi dan Sosial Universitas Jenderal Achmad Yani Yogyakarta

Jln. Siliwangi Ringroad Barat, Yogyakarta 55293 INDONESIA

Email: ^{1*}alfirnarizqi@gmail.com, ²adliazhafarina@gmail.com, ³sapnanda9@gmail.com, ⁴neethajhary@gmail.com

Abstract— *The APJII 2022 survey noted that 98.02% of the internet is used by generation Z young people to access social media. The existence of social media certainly has a social impact, one of its as a place of bullying. Problems arise in the process of identifying victims and perpetrators of cyberbullying, besides difficulties in the meaning of words that contain elements of bullying can cause multiple interpretations that affect the investigation process. To overcome these problems, an initial identification effort is needed by grouping topics of conversation using the text mining method. This method can be started by preprocessing techniques, because the conversation data is in the text form. Next, topic grouping is carried out using the K-means clustering algorithm. Based on the results of the calculation, comment groups containing bullying are then mapped with Rapidminer tool. The contribution of this research is as the first step in the implementation of cybersecurity at the stage of identifying digital evidence for investigation purposes. This research resulted in mapping in the form of clusters. The results of the analysis on 2 clusters showed a grouping pattern in certain terms that were consistently grouped and contained the same term, namely the term "Nanyi", so that it can be concluded early that the topic being discussed in the Video is about singing activities or elements of conversation in the context of singing. While in other clusters, the terms in groups are more diverse because of the uniqueness of each term.*

Intisari— Survey APJII 2022 mencatat 98,02% internet digunakan oleh kalangan anak muda generasi Z untuk mengakses media sosial. Keberadaan media sosial tentu saja memiliki dampak sosial, salah satunya sebagai tempat perundungan atau *bullying*. Permasalahan muncul dalam proses identifikasi korban dan pelaku *cyberbullying*, selain itu kesulitan dalam makna dan maksud kata yang mengandung unsur *bullying* dapat menimbulkan multitafsir yang mempengaruhi dalam proses investigasi. Untuk mengatasi permasalahan tersebut, diperlukan sebuah upaya identifikasi awal dengan mengelompokkan topik pembicaraan menggunakan metode *text mining*. Cara ini dapat dimulai dengan melakukan teknik *preprocessing*, dikarenakan data percakapan berbentuk teks. Selanjutnya dilakukan pengelompokkan topik menggunakan algoritma *K-means clustering*. Berdasarkan hasil perhitungan kelompok komentar yang mengandung perundungan kemudian dipetakan dengan bantuan *tools Rapidminer*. Kontribusi dari penelitian ini adalah sebagai langkah awal implementasi keamanan siber pada tahap identifikasi bukti digital untuk keperluan investigasi dan persidangan. Penelitian ini menghasilkan pemetaan dalam bentuk kluster. Hasil analisis pada 2 kluster terlihat pola pengelompokkan pada *term* tertentu yang secara konsisten terkelompok dan berisi *term* yang sama yaitu *term* "Nanyi", sehingga dapat ditarik

kesimpulan awal bahwa topik yang sedang menjadi pembahasan dalam Video adalah seputar kegiatan bernyanyi atau unsur perbincangan dalam konteks bernyanyi. Sedangkan pada kluster lain *term* yang berkelompok lebih beragam karena keunikan masing-masing *term*.

Kata Kunci— *Cyberbullying, Clustering, K-Means, Topic Modelling, Youtube*

I. PENDAHULUAN

Internet menjadi infrastruktur utama menuju dunia siber. Survey APJII 2022 mencatat 98,02% internet digunakan untuk mengakses media sosial dan penggunaannya didominasi oleh kalangan generasi Z, dimana media yang paling sering dikunjungi adalah Youtube sebesar 63,02%[1]. Media sosial adalah bentuk kemajuan teknologi yang mempengaruhi masyarakat dalam aspek gaya hidup, cara pandang dan budaya[2]. Di kalangan anak muda generasi Z media sosial yang populer seperti Tiktok, Instagram, YouTube.

Media sosial yang paling banyak dikunjungi salah satunya adalah YouTube[1]. Layanan yang diberikan YouTube sangat beragam dengan interaksi komentar dan *emoticon* sehingga pengguna dapat mengekspresikan dirinya secara interaktif. Fitur komentar dapat ditemukan di video, YouTube Live, dan YouTube Shorts. Keberadaan media sosial tentu saja memiliki dampak positif dan negatif, sebagai salah satu dampak negatif dari keberadaan media sosial yaitu menggunakannya sebagai tempat perundungan atau *bullying*.

Bullying merupakan perilaku psikologis yang menjadi kejahatan verbal dan lebih menyakitkan dibandingkan kekerasan fisik seperti melecehkan, menghina, mengancam, merendahkan, atau menyakiti dan berujung pada penghinaan terhadap orang lain[3]. Perundungan yang dilakukan dalam aktivitas siber seperti komentar negatif pada postingan, pesan personal yang tak baik, serta mengolok-olok yang aktivitas tersebut terjadi di dunia siber disebut sebagai *cyberbullying*.

Merujuk pada undang-undang nomor 11 Tahun 2008 tentang Informasi dan Transaksi Elektronik (UU ITE). Tindakan menunjukkan penghinaan terhadap orang lain tercermin pada pasal 27 ayat (3) UU ITE yang berbunyi: "Setiap orang dengan sengaja dan tanpa hak mendistribusikan dan/atau mentransmisikan dan/atau membuat dapat diaksesnya informasi elektronik dan/atau dokumen elektronik yang memiliki muatan penghinaan dan/atau pencemaran nama baik".

Permasalahan muncul dalam proses identifikasi korban dan pelaku *cyberbullying*, selain itu kesulitan dalam makna dan maksud kata yang mengandung unsur *bullying* dapat menimbulkan multitafsir yang mempengaruhi dalam proses investigasi. Untuk mengatasi permasalahan tersebut, dibutuhkan keahlian di bidang forensika digital, tujuannya untuk mendukung langkah investigasi dan pencarian barang bukti pada kejahatan *cyberbullying*.

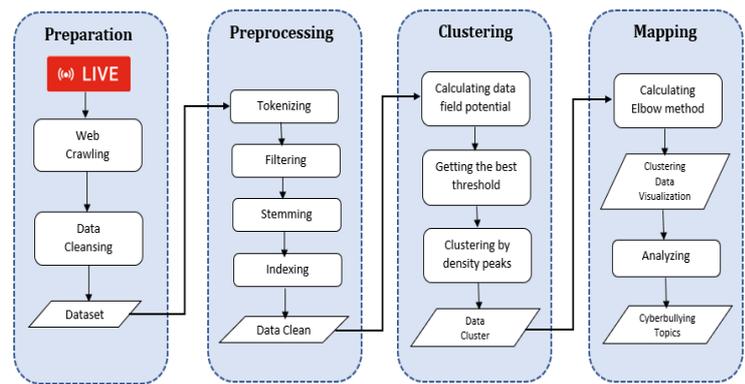
Sejumlah penelitian yang berkaitan dengan media sosial dan *cybersecurity*, dilakukan pada data Twitter[4] dan Youtube[5]. Maraknya kasus kejahatan melalui *platform* media sosial ditindaklanjuti secara khusus menginvestigasi kelompok kata sebagai alat bukti tindak kejahatan verbal[6]. Tindak kejahatan verbal *bullying* dapat mengancam pelakunya sebagai pertanggungjawaban pidana[7]. Sebagian besar analisis kata yang mengandung *cyberbullying* dalam bentuk teks, sehingga diolah menggunakan *textmining*, seperti penelitian[8] yang mengaplikasikan algoritma *K-means clustering* menggunakan *tools* Rapidminer. Penelitian dengan *textmining* untuk menganalisis pengaruh sebuah hastag di media sosial[9]. *Textminig* juga digunakan untuk pengelompokan topik menggunakan algoritma *clustering*[10], pengelompokan data dapat menjadi kontribusi dalam mengetahui topik[11]. Algoritma *clustering* dapat dikombinasikan dengan *cosine similarity*[12], maupun TF-IDF[13].

K-means merupakan metode yang bekerja dengan cara membagi himpunan data ke dalam kelompok, penentuan kelompok dilihat dari jarak terdekat antara data dengan *centroid*. Penelitian dengan algoritma *K-Means clustering* sebagai cara identifikasi awal *cyberbullying*[14],[3],[15],[16],[17]. Untuk mengukur kualitas dan optimalisasi *cluster* digunakan Elbow method[18],[19],[20],[21]. Berdasarkan sejumlah studi literatur yang telah dibahas, media sosial merupakan objek yang populer untuk dieksplorasi, khususnya komentar di Youtube. Selain memiliki fitur komentar secara *asynchronous*, YouTube juga dilengkapi dengan fitur komentar secara langsung (*live*). Tidak jarang pada komentar ditemukan percakapan yang tidak sesuai bahkan menggunakan kata-kata yang kurang pantas.

Pada penelitian ini, agar dapat mendeteksi kejahatan *cyberbullying* dalam sebuah percakapan pada komentar di *live* YouTube, diperlukan sebuah upaya identifikasi awal dengan mengelompokkan topik pembicaraan menggunakan metode *text mining*. Dikarenakan data percakapan berbentuk teks, cara ini dapat dimulai dengan melakukan teknik *preprocessing*. Selanjutnya dilakukan pengelompokan topik berdasarkan kata yang sudah teridentifikasi menggunakan algoritma *K-means clustering*. Berdasarkan hasil perhitungan kelompok komentar yang mengandung perundungan kemudian dipetakan atau dimodelkan dengan bantuan *tools* Rapidminer. Kontribusi dari penelitian ini adalah sebagai langkah awal implementasi keamanan siber pada tahap identifikasi bukti digital untuk keperluan investigasi dan persidangan.

II. METODE

Penelitian ini menggunakan data berbasis teks yang dihimpun berasal dari komentar *live* pada sebuah kanal Youtube. Adapun tahap dalam penelitian ini seperti yang ditunjukkan pada Gambar 1.



Gambar. 1 Tahap Penelitian

Berdasarkan Gambar 1, tahapan dalam penelitian ini meliputi 4 tahap sebagai berikut:

A. Tahap Persiapan

Pada tahap ini telah dilakukan studi literatur untuk kajian awal penelitian. Selanjutnya akan dilakukan penentuan sumber data berupa channel Youtube untuk dilakukan tahap pengumpulan Data. Data yang dihimpun merupakan data komentar berbasis teks. Proses menghimpun data menggunakan teknik *crawling*, yaitu proses pemindaian sebuah konten pada *website*. Data yang telah didapatkan akan filter terlebih dahulu dari *noise* untuk siap menjadi dataset dalam penelitian.

B. Tahap Pengolahan Data (Preprocessing)

Data dalam bentuk teks yang telah dihimpun akan melalui beberapa proses, tujuan proses ini adalah ekstraksi data dengan *preprocessing* yang meliputi *Tokenizing*, *Filtering*, *Stemming* dan *Indexing*. Hasil akhir *preprocessing* berupa data bersih yang siap diolah[22].

C. Tahap Analisis Data (Clustering)

Data yang telah diekstraksi kemudian dianalisis dengan metode *clustering* menggunakan algoritma *K-means clustering* sebagai formula dalam menghitung kelompok kata. *Clustering* merupakan teknik *data mining* yang mengelompokkan data berdasarkan pada objek/karakteristik yang sama kemudian dikumpulkan dalam satu kluster dan kluster tersebut berbeda dengan kluster lainnya. *K-Means* digunakan untuk mengelompokkan karakter berdasarkan karakter yang telah didefinisikan. Persamaan *K-Means* seperti yang ditunjukkan pada persamaan 1:[9]

$$d(X_j, C_j) = \sqrt{\sum_{j=1}^n (X_j - C_j)^2} \quad (1)$$

Keterangan:

d = distance

n = number of objects

j = start from 1 to n

X_j = feature object to j with respect to x

C_j = centroid feature to j

D. Tahap Pemetaan Hasil (Mapping)

Hasil klusterisasi akan dipetakan dengan bantuan *tools* Rapidminer dalam bentuk visual *scatter* diagram maupun *wordcloud*.

III. HASIL DAN PEMBAHASAN

Hasil penelitian ini berupa pemetaan kelompok kata sebagai identifikasi awal *cyberbullying*. Adapun hasil dari masing-masing tahapan sebagai berikut:

A. Tahap Persiapan

Pada penelitian ini, pengumpulan data berasal dari komentar pada sebuah YouTube yang merupakan data berbasis teks. Pengumpulan data dilakukan dengan teknik *crawling* yang bertujuan mengambil data tertentu, adapun data komentar berasal dari channel berikut: <https://www.youtube.com/shorts/RABxiBvgxpU>. Hasil *crawling* diperoleh 100 data yang berisi komentar yang akan menjadi sumber data untuk diolah dalam ranah *text mining*.

B. Tahap Pengolahan Data (Preprocessing)

Hasil *crawling* disimpan dalam bentuk *.csv* untuk selanjutnya diolah dengan teknik *preprocessing* yang terdiri dari *case folding*, *tokenizing*, *stopword removal*, *stemming* dan *indexing*. Proses *preprocessing* menggunakan bahasa Python. Berikut tampilan *preprocessing* yang ditunjukkan pada Gambar 2.

```
[ ] import pandas as pd
import numpy as np

tweet_data = pd.read_csv("data.csv")

tweet_data.head()

      tweet  Unnamed: 1
0  Ini sih masalah teknik. Berarti tehnik bernyan...  view comment
1  Sebenarnya suara nih orang tuh bagus banget....  view comment
2  Dasar yang edit dan buat cerita pintar banget!...  view comment
3  Jam terbang GK pernah bohong ..  view comment
4  Nada tinggi itu susah  view comment

tweet_data['tweet'] = tweet_data['tweet'].str.lower()

print('Case Folding result : \n')
print(tweet_data['tweet'].head(5))
print('\n\n')
```

Gambar. 2 Tahap *Preprocessing*

C. Tahap Analisis Data (Clustering)

Data yang telah *dipreprocessing* kemudian disimpan dalam format *.csv*. Dikluster menggunakan *tools* Rapidminer dengan pengujian 2 kluster, 3 kluster, 5 kluster, 7 kluster dan 10 kluster. Berikut adalah cuplikan masing-masing perhitungan pada Tabel 1:

Tabel. 1 Hasil Kluster

No	Kluster	Total Item	Keterangan
1	2	488	Cluster 0: 468 items Cluster 1: 20 items

No	Kluster	Total Item	Keterangan
2	3	488	Cluster 0: 467 items Cluster 1: 20 items Cluster 2: 1 items
3	5	488	Cluster 0: 464 items Cluster 1: 20 items Cluster 2: 1 items Cluster 3: 1 items Cluster 4: 2 items
4	7	488	Cluster 0: 456 items Cluster 1: 20 items Cluster 2: 3 items Cluster 3: 1 items Cluster 4: 6 items Cluster 5: 1 items Cluster 6: 1 items
5	10	488	Cluster 0: 422 items Cluster 1: 15 items Cluster 2: 16 items Cluster 3: 20 items Cluster 4: 6 items Cluster 5: 1 items Cluster 6: 2 items Cluster 7: 5 items Cluster 8: 4 items Cluster 9: 2 items Cluster 10: 1 items

D. Tahap Pemetaan Hasil (Mapping)

Hasil klustering yang dilakukan bertujuan untuk menemukan kelompok kata yang memiliki kesamaan makna atau konteks dikelompokkan dalam kluster yang sama. Setiap kelompok kata yang terbentuk oleh algoritma akan diwakili oleh pusat kluster (*centroid*) dan berisi sejumlah kata. Kata-kata yang paling dekat dengan pusat kluster adalah yang paling mendefinisikan kluster tersebut. Menemukan pola dan hubungan antara kata-kata dalam teks. Kesamaan makna kata dianalisis dari *term* yang paling sering muncul bersama dalam teks, sehingga apabila *term* tersebut memiliki makna yang sama, maka akan terkelompok dalam kluster yang sama. Dalam konteks ini dapat diidentifikasi tema atau topik tertentu dan mengidentifikasi pola sentimen.

Berikut adalah hasil pemetaan *term* yang divisualisasikan dalam bentuk diagram *scatter* menggunakan Rapidminer yang ditunjukkan pada Gambar 3.

- center.org/index.php/snrakt2017/article/download/10/9.
- [9] M. Habibi and P. W. Cahyo, "Clustering User Characteristics Based on the influence of Hashtags on the Instagram Platform," *IJCCS (Indonesian J. Comput. Cybern. Syst.*, vol. 13, no. 4, p. 399, 2019, doi: 10.22146/ijccs.50574.
- [10] H. Irsyad and M. R. Pribadi, "Implementasi Text Mining Dalam Pengelompokan Data Tweet Pertanian Indonesia Dengan K-Means," *Kurawal - J. Teknol. Inf. dan Ind.*, vol. 3, no. 2, pp. 164–172, 2020, doi: 10.33479/kurawal.v3i2.347.
- [11] F. Kolini and L. Janczewski, "Clustering and topic modelling: A new approach for analysis of national cybersecurity strategies," *Proc. of 21st Pacific Asia Conf. Inf. Syst. "Societal Transform. Through IS/IT"*, PACIS 2017, 2017.
- [12] N. W. Utami and I. G. J. Eka Putra, "Text Minig Clustering Untuk Pengelompokan Topik Dokumen Penelitian Menggunakan Algoritma K-Means Dengan Cosine Similarity," *J. Inform. Teknol. dan Sains*, vol. 4, no. 3, pp. 255–259, 2022, doi: 10.51401/jinteks.v4i3.1907.
- [13] R. Khan, Y. Qian, and S. Naeem, "Extractive based Text Summarization Using KMeans and TF-IDF," *Int. J. Inf. Eng. Electron. Bus.*, vol. 11, no. 3, pp. 33–44, 2019, doi: 10.5815/ijeeb.2019.03.05.
- [14] K. P. Sinaga and M. S. Yang, "Unsupervised K-means clustering algorithm," *IEEE Access*, vol. 8, pp. 80716–80727, 2020, doi: 10.1109/ACCESS.2020.2988796.
- [15] Aderibigbe, "Pemodelan Deteksi Cyber Bullying Pada Jejaring Sosial Twitter," *Energies*, vol. 6, no. 1, pp. 1–8, 2018, [Online]. Available: <http://journals.sagepub.com/doi/10.1177/1120700020921110%0Ahttps://doi.org/10.1016/j.reuma.2018.06.001%0Ahttps://doi.org/10.1016/j.arth.2018.03.044%0Ahttps://reader.elsevier.com/reader/sd/pii/S1063458420300078?token=C039B8B13922A2079230DC9AF1A333E295FCD8>.
- [16] M. Ahmed, R. Seraj, and S. M. S. Islam, "The k-means algorithm: A comprehensive survey and performance evaluation," *Electron.*, vol. 9, no. 8, pp. 1–12, 2020, doi: 10.3390/electronics9081295.
- [17] N. F. Hasan, "Deteksi Cyberbullying pada Facebook Menggunakan Algoritma K-Nearest Neighbor," *J. Smart Syst.*, vol. 1, no. 1, pp. 35–44, 2021, doi: 10.36728/jss.v1i1.1605.
- [18] H. Humaira and R. Rasyidah, "Determining The Appropriate Cluster Number Using Elbow Method for K-Means Algorithm," 2020, doi: 10.4108/eai.24-1-2018.2292388.
- [19] M. A. Syakur, B. K. Khotimah, E. M. S. Rochman, and B. D. Satoto, "Integration K-Means Clustering Method and Elbow Method for Identification of1. Syakur MA, Khotimah BK, Rochman EMS, Satoto BD. Integration K-Means Clustering Method and Elbow Method for Identification of the Best Customer Profile Cluster. IOP Conf Ser Mat," *IOP Conf. Ser. Mater. Sci. Eng.*, vol. 336, no. 1, 2018, doi: 10.1088/1757-899X/336/1/012017.
- [20] M. Cui, "on the Elbow Method," pp. 5–8, 2020, doi: 10.23977/accaf.2020.010102.
- [21] J. D'Silva and U. Sharma, "Unsupervised Automatic Text Summarization of Konkani Texts using K-means with Elbow Method," *Int. J. Eng. Res. Technol.*, vol. 13, no. 9, pp. 2380–2384, 2020, doi: 10.37624/ijert/13.9.2020.2380-2384.
- [22] A. R. Lahitani, "Automated Essay Scoring menggunakan Cosine Similarity pada Penilaian Esai Multi Soal," *J. Kaji. Ilm.*, vol. 22, no. 2, pp. 107–118, 2022, doi: 10.31599/jki.v22i2.1121.