

# Identifikasi Topik Hangat di Media Berita Menggunakan Latent Dirichlet Allocation

Aulia Cisatra<sup>1</sup>, Daffa Daris Mahendra Ansori<sup>2</sup>, Dara Nasywa Fathya Afiqah Akmal<sup>3</sup>, Slamet Ramadhani<sup>4</sup>,  
Nur Aini Rakhmawati<sup>5</sup>

<sup>1,2,3,4,5</sup> Program studi S-1 Sistem Informasi, Institut Teknologi Sepuluh Nopember

<sup>1</sup>[cisatra.205026@mhs.its.ac.id](mailto:cisatra.205026@mhs.its.ac.id)

<sup>2</sup>[daffa.205026@mhs.its.ac.id](mailto:daffa.205026@mhs.its.ac.id)

<sup>3</sup>[dara.2050261@mhs.its.ac.id](mailto:dara.2050261@mhs.its.ac.id)

<sup>4</sup>[5998231017@student.its.ac.id](mailto:5998231017@student.its.ac.id)

<sup>5</sup>[nur.aini@is.its.ac.id](mailto:nur.aini@is.its.ac.id)

**Abstrak**— Akses informasi di era digital saat ini sebagian besar bersumber dari media daring. Namun, identifikasi topik aktual yang tengah hangat dibicarakan publik kian menjadi tantangan untuk menentukan topik berita yang relevan dan menarik perhatian pembaca. Pada studi ini, dikembangkan sistem cerdas berbasis pembelajaran mesin guna membantu portal berita *online* dalam memahami minat publik terhadap beragam topik. Memanfaatkan pendekatan *natural language processing*, model *latent dirichlet allocation* akan digunakan untuk analisis dan klasifikasi berita berdasarkan topiknya. Data teks digital yang digunakan berasal dari beragam portal berita nasional melalui akses *application programming interface* (API). Melalui *text mining* dan *machine learning*, ekstraksi topik utama secara otomatis dari data teks yang besar dan tidak terstruktur dapat dikenali secara cepat juga akurat. Metode ini memungkinkan pengenalan atas topik-topik yang mendominasi berita daring, sehingga memungkinkan pembaca, peneliti, hingga praktisi media untuk tetap terhubung dengan isu terkini. Hasil studi menunjukkan tiga topik terpopuler berdasarkan dominasi token secara berurutan adalah KPK (36.6%), cawapres (32.8%), dan jakarta (30.6%). Penelitian ini diharapkan dapat meningkatkan daya saing portal berita dalam menyajikan konten aktual sesuai preferensi pembaca. Selain itu, hasil studi dapat menjadi acuan bagi peneliti lain dalam bidang serupa di masa mendatang.

**Kata Kunci**— Latent Dirichlet Allocation, Media Berita, Natural Language Processing, Topic Modelling, Topik Hangat.

## I. PENDAHULUAN

Merupakan hak setiap manusia atas tersedianya akses ke berita terbaru dari seluruh penjuru dunia. Terlebih di era digital, siaran berita melalui media dalam jaringan pun tidak menjadi pengecualian melainkan salah satu sumber utama dalam persebaran informasi [1]. Fasilitas portal berita daring memungkinkan masyarakat untuk terus terhubung dan mendapatkan pembaruan informasi dalam berbagai bidang termasuk sektor politik, ekonomi, dan budaya secara mudah. Melalui media berita daring, pembaca dapat turut berpartisipasi aktif dalam forum diskusi dan mengutarakan gagasan terkait topik berita melalui kolom komentar dan media sosial. Meskipun diskusi ini dapat saling mencerdaskan antar individu akibat pertukaran perspektif, tantangan dalam mengidentifikasi topik yang mendominasi perbincangan publik menjadi tantangan tersendiri akibat banyaknya informasi yang tersedia. Oleh karena itu, dibutuhkan sistem

yang dapat secara cerdas menganalisis data teks digital guna mengetahui topik perbincangan utama.

Berbagai pendekatan telah dikembangkan sebagai upaya dalam mengatasi kendala indentifikasi topik hangat ini. Salah satunya adalah melalui analisis sentimen dan *topic modelling* berbasis pemrosesan bahasa alami atau *natural language processing* [2]. Meskipun demikian, perkembangan dalam siklus pemberitaan serta dinamika minat masyarakat turut menjadi tantangan dalam penelitian sehingga media berita pun memerlukan pendekatan yang lebih cermat dan adaptif terhadap perkembangan terkini. Guna mengatasi tantangan tersebut, penelitian ini mengembangkan sistem indentifikasi dan klasifikasi topik berita demi mendorong hadirnya sajian informasi intelektual dan berkualitas sesuai minat pembaca.

Pada studi ini dikembangkan sistem cerdas berbasis *natural language processing* (NLP) menggunakan metode *topic modelling* dan pemanfaatan algoritma *latent dirichlet allocation* (LDA) sebagai solusi. Bertujuan untuk merepresentasikan preferensi publik secara nyata, data teks digital bersumber dari kumpulan informasi berita terkini oleh portal berita daring di Indonesia. Dengan memanfaatkan *application programming interface* (API) sebagai sumber data, sistem ini dirancang untuk dapat bekerja secara otomatis dan adaptif dalam analisis sejumlah besar data teks guna ekstraksi daftar topik paling banyak diperbincangkan publik [3].

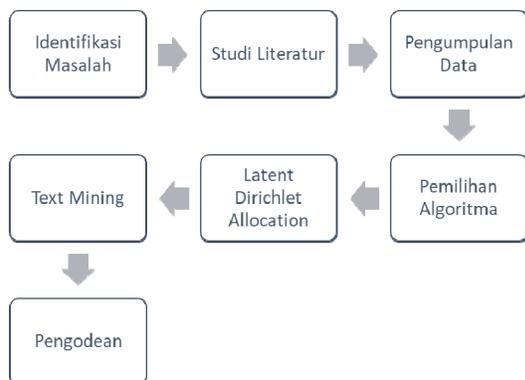
Penelitian serupa terkait identifikasi topik pada portal media daring pernah dikembangkan pada 2022 lalu [4]. Menggunakan *topic modelling* dengan algoritma LDA, penelitian tersebut mengumpulkan data berita dari salah satu akun Twitter milik media berita Detik.com dalam kurun waktu satu tahun. Adapun studi tersebut memiliki potensi pengembangan dengan penggunaan data yang hanya bersumber dari satu portal berita dan rentang waktu yang terbatas. Sehingga pada penelitian ini diupayakan perluasan sumber data yang digunakan dengan tetap terbatas pada situs portal berita dalam negeri oleh situs antaranews.com, cncbindonesia.com, cnnindonesia.com, jpnn.com, merdeka.com, dan beberapa situs lainnya. Selain itu, penelitian ini pun akan menggunakan rentang serta kategori topik pembicaraan yang lebih luas demi relevansi atas kondisi masyarakat saat ini.

Algoritma LDA ini pun pernah digunakan dan diteliti pada 2022 lalu untuk kasus analisis topik pembahasan dominan melalui pemantauan tagar covidindonesia [5]. Penelitian tersebut dilakukan atas analisis hasil *text mining* 84 takarir (*caption*) dari berbagai pengguna sosial media Instagram. Dalam menentukan jumlah topik yang optimal, penelitian tersebut menggunakan penilaian berdasarkan *perplexity* dan *topic coherence* yang menunjukkan lima topik teratas suatu unggahan video. Identifikasi topik atas konten tersebut antara lain meliputi covidindonesia, covid\_19, pandemi di Indonesia, dan pembahasan mutasi virus covid-19.

Selain itu, penelitian lainnya terkait identifikasi konten oleh *virtual YouTuber* (VTuber) terkenal turut menerapkan algoritma LDA dengan mempertimbangkan nilai *perplexity* dan *topic coherence* dalam validasi penentuan topik [6]. Analisis terhadap 4312 video dari 10 kanal VTuber teratas di Indonesia tersebut menghasilkan lima topik tayangan teratas para VTuber Indonesia berupa gim Minecraft dan *reading donation*, gim Apex Legend disertai kolaborasi dengan VTuber lain, siaran langsung video gim, *cover* lagu, serta tayangan gim *multiplayer* seperti Raft atau Phasmophobia.

Pemanfaatan *topic modelling* pada jejaring sosial seperti Twitter ataupun portal berita tentu tidak terhitung jumlahnya. Adanya kumpulan penelitian terdahulu menggunakan berbagai pendekatan terhadap area penelitian yang bervariasi pun mendukung gagasan bahwa kombinasi penggunaan *topic modelling* dan algoritma LDA dapat menghasilkan akurasi yang sama baiknya.

Pada penelitian ini, pemanfaatan *topic modelling* dalam identifikasi topik hangat pada berbagai portal berita daring serta penggunaan algoritma LDA yang mampu mengelompokkan, menghubungkan, dan memproses data berita dalam jumlah besar akan digunakan. Selain itu, dipastikan pula bahwa seluruh data berada dalam kategori terbaru demi hasil analisis topik yang relevan dengan masanya. Adapun studi ini diharapkan dapat meningkatkan daya saing media *online* serta memberikan kontribusi dalam pemantauan tren juga isu penting yang tengah berlangsung di Indonesia. Selain itu, diharapkan juga penelitian ini dapat menjadi rujukan dan inspirasi bagi para praktisi akademik untuk penelitian yang berkaitan.



Gbr 1. Metodologi Penelitian

## II. METODE PENELITIAN

Pada penelitian ini digunakan pembelajaran mesin, tepatnya pendekatan *natural language processing* (NLP). Metodologi studi ini antara lain meliputi tahapan identifikasi masalah, pembelajaran literatur, pengumpulan data, pemilihan algoritma, penggunaan *latent dirichlet allocation* (LDA), hingga *text mining* dan pengodean sebagaimana [Gambar 1].

### A. Identifikasi Masalah

Tahapan identifikasi masalah dilakukan guna untuk mengetahui kendala yang melatarbelakangi penelitian, yaitu seiring dengan derasnya peredaran informasi di portal berita, seringkali terjadi kekeliruan dalam penafsiran suatu informasi dan intisari berita dari media sosial. Sehingga pada penelitian ini diajukan pengembangan sistem identifikasi topik hangat menggunakan *natural language processing*.

### B. Studi Literatur

Dalam bagian ini akan diuraikan kajian pustaka yang relevan dengan penelitian identifikasi topik hangat di media berita *online*. Kajian ini membahas beberapa konsep kunci yang mendukung landasan teori dan metodologi penelitian.

#### 1) Data preprocessing

Pra-pemrosesan data merupakan tahapan persiapan untuk menghasilkan kualitas pengelompokan (*clustering*) yang baik [7]. Rangkaian tahapan yang dilakukan pada penelitian ini adalah (1) mengubah kalimat menjadi kata; (2) menghilangkan beberapa kata dalam bahasa Indonesia dan bahasa Inggris yang tidak memiliki arti, seperti kata *'using'*, *'of'*, *'the'*, *'in'*, *'on'*, *'as'*, *'and'*, dan *'based'* pada bahasa Inggris dan kata *'berbasis'*, *'kasus'*, *'dan'*, juga *'pada'* dalam bahasa Indonesia [12]; (3) mengubah susunan kalimat dalam bentuk bigram, misalnya konversi kalimat "Inovasi Media Pembelajaran Sain Teknologi di SMP Berbasis Mikrokontroler" menjadi [inovasi, media], [pembelajaran, sain], [teknologi, di], [SMP, berbasis], dan juga [mikrokontroler] [8].

#### 2) Natural language processing (NLP)

Pemrosesan bahasa alami atau NLP adalah cabang keilmuan dari teknologi kecerdasan buatan yang berfokus pada pemahaman dan pengolahan bahasa manusia secara alami oleh komputer [2]. Keilmuan ini mencakup sejumlah teknik untuk analisis dan pemrosesan data berbasis teks, seperti tokenisasi, entitas nama, analisa sentimen, hingga pemodelan topik. Dalam konteks identifikasi topik berita, NLP dapat digunakan untuk mengelompokkan artikel berita berdasarkan topik bahasanya.

#### 3) Topic modelling

Pemodelan topik atau *topic modelling* memungkinkan identifikasi kategori dalam koleksi dokumen teks dengan pengelompokan (*clustering*) berdasarkan kemunculan kata kunci untuk setiap topik. Salah satu metode paling populer dalam *topic*



