



Degree Program in Data Science and Management

Course of Machine Learning

"Automating News":
The Role of AI in
Journalism

Prof. Giuseppe Italiano

SUPERVISOR

Prof. Irene Finocchi

CO-SUPERVISOR

ID 782571
Navarra Filippo

CANDIDATE

Academic Year 2024/2025

Contents

1	Introduction	7
1.1	Background and Context	7
1.2	Research Aim and Objectives	8
1.3	Research Aim and Objectives	8
1.4	Research Questions	9
1.5	Methodological Overview	9
1.6	Structure of the Thesis	10
2	Journalism in the Digital Era	11
2.1	Defining Automated Journalism	11
2.1.1	Conceptual Foundations	11
2.1.2	Historical Evolution and Trends	12
2.2	Technological Foundations	13
2.2.1	Natural Language Processing (NLP): Extraction and Summarization	13
2.2.2	Natural Language Generation (NLG): Report Writing and Rewording	14
2.2.3	Machine Learning Techniques: Classification, Clustering, Predictive Models	15
2.2.4	Web Scraping, Data Ingestion, and Real-Time Feeds	15
2.3	Benefits of AI in Newsroom Workflows	16
2.3.1	Efficiency and Speed	16
2.3.2	Scalability and Personalization	16
2.3.3	Innovative Content Forms	17
2.4	Risks and Challenges	17
2.4.1	Data Bias and Algorithmic Fairness	17
2.4.2	Transparency and Editorial Accountability	18
2.4.3	Impact on Labor and Editorial Roles	19
2.4.4	Content Farms	19
2.5	Summary of Literature Gaps	21
3	The Problem: Why Agencies Are Turning to AI	23
3.1	Changing Information Ecosystem	23
3.1.1	24/7 News Cycle and Content Saturation	23
3.1.2	Rise of Data-Driven and Personalized Journalism	24
3.2	Bottlenecks in Traditional News Workflows	26
3.2.1	Manual Curation and Time Constraints	26

3.2.2	Limitations in Real-Time Reporting	26
3.3	Operational Pressures	27
3.3.1	Budget Constraints and Staff Reductions	27
3.3.2	Competition with Non-traditional Publishers	27
3.4	Editorial Challenges	29
3.4.1	Quality Control at Scale	29
3.4.2	Maintaining Accuracy under Pressure	30
3.5	Summary: Need for Workflow and Information Flow Automation	31
4	The Solution: Applying AI and GenAI to News Automation	33
4.1	Overview of AI Solutions in the News Industry	33
4.1.1	Workflow Automation with AI Agents	33
4.1.2	Integration of Data Pipelines and AI Systems	34
4.2	Generative AI for News Content	35
4.2.1	Text Generation and Rewriting Models	35
4.2.2	Automated Headline and Summary Generation	37
4.3	Machine Learning for Editorial Assistance	39
4.3.1	Classification of News Content by Topic and Tone	39
4.3.2	Predictive Analytics for Newsworthiness and Audience Engagement	40
4.4	Human-in-the-Loop vs. Full Automation	40
5	Designing the Journalist Journey: Tracing a News Article in an AI-Powered Newsroom	42
5.1	Introduction	42
5.2	News Gathering	42
5.3	News Processing	47
5.3.1	Named Entity Recognition	49
5.3.2	Sentiment Analysis	49
5.3.3	Topic Modelling	50
5.4	News Generation	51
5.4.1	Generated Article	53
5.4.2	Human-in-the-Loop Validation	53
5.5	News Distribution	54
5.5.1	Publishing	54
5.5.2	Translating & Targeting	54
6	Discussions and Conclusions	56
6.1	Revisiting Research Questions	56
6.2	Ethical and Regulatory Considerations	57
6.2.1	Bias, Transparency, and Trust	57
6.2.2	The European Legislative Landscape (AI Act, GDPR, etc.)	58
6.3	Final Thoughts and Conclusion	59

List of Figures

2.1	Today's Quakebot online published content - LA Times	12
2.2	NewsRooms usage of AI Tools	14
2.3	Stakeholders and GenAI	20
2.4	Number of AI generated News Sites by NewsGuard	21
3.1	Reasons why adopting AI in newsrooms	24
3.2	The workflow of news recommendation system	25
3.3	Newsrooms employment in USA - 2008-2020	28
3.4	An example of Non-traditional competitor's content	29
3.5	Media Consumption Trends	30
3.6	Attitude towards AI produced news	31
3.7	AI active strategies in newsrooms	32
4.1	BERT Model Architecture	36
5.1	Time Series Plot for #EstaniaProtest	45
5.2	Sentiment Analysis from social media posts	52

List of Tables

5.1	MTEB (Multilingual) Benchmark Results	48
5.2	Performance comparison of transformer-based models on the SST-2 dataset.	50

Listings

- 5.1 Example of raw XML feed from Estania 43
- 5.2 Trending Posts from r/politics 44
- 5.3 Simulated X XML response for #EstaniaProtest 45
- 5.4 Simulated Instagram XML response for #EstaniaProtest 46

Chapter 1

Introduction

1.1 Background and Context

Journalism has undergone radical transformations, shaped by successive waves of technological innovation from the invention of the printing press to the rise of broadcast and, more recently, the explosion of digital platforms. In the current phase, artificial intelligence (AI) is emerging as a transformative force, introducing new opportunities for content creation, real-time analysis, and scalable distribution. This phase, often referred to as *Automated Journalism*, involves the use of algorithmic systems to produce news stories with minimal human intervention. As technology advanced into the digital era, the internet changed the way news was created, consumed, and disseminated, leading to the emergence of online journalism. Currently, the incorporation of artificial intelligence and automation in newsrooms is transforming the environment once more, facilitating quicker, data-informed news generation. From print to broadcast to digital and automated journalism, the industry has consistently evolved with technological progress, transforming how stories are narrated and experienced.

In this thesis, we will refer to this phenomenon as *Automated Journalism*, a term that captures the growing role of artificial intelligence in the news production process. While multiple definitions exist in both academic and industry discourse - some describing a more limited degree of AI involvement, such as *Algorithmic Journalism* (K. N. Dörr, 2015), and others encompassing broader applications, like *Robot Journalism* - this work will focus on the automation of news generation through artificial intelligence. By utilizing computational tools to analyse data and produce narratives, automated journalism is reshaping the media landscape, raising critical questions about the balance between efficiency and editorial oversight.

While AI has already found limited roles in newsrooms (e.g., through Natural Language Generation for financial reporting), recent developments in generative AI (GenAI), large language models (LLMs), and machine learning pipelines are enabling more comprehensive automation across the news production chain. These developments not only accelerate content production but also pose significant editorial, ethical, and regulatory challenges. For example, the increased use of black-box models raises concerns about transparency, accountability, and the preservation of journalistic integrity.

This thesis addresses a key problem: **How can AI systems be responsibly and effectively integrated into newsroom workflows without compromising core editorial values?** The relevance of this question stems from an ongoing shift in journalistic labor, where AI systems are no longer isolated tools but are becoming embedded across sourcing, drafting, publishing, and audience analytics.

Although existing literature has analyzed components of this shift from the early stages of algorithmic news writing (e.g., K. N. Dörr, 2015; Graefe, 2016) to discussions about bias and automation risks, most studies treat these issues in isolation. There is a lack of comprehensive models that bridge technological implementation with editorial practice, particularly under the pressures of real-world newsroom constraints and regulatory frameworks such as the EU AI Act or GDPR.

This thesis advances the state of the art by presenting an applied, systems-oriented exploration of AI in journalism, culminating in a simulated end-to-end newsroom workflow that illustrates how automation can be balanced with human oversight. Through both technical analysis and ethical reflection, it contributes a holistic model for thinking about AI in journalism as an integrated sociotechnical system.

1.2 Research Aim and Objectives

The aim of this research is to critically examine and model the integration of artificial intelligence systems—particularly generative and assistive AI—into journalistic workflows. This includes a dual focus: (1) understanding the technical underpinnings of automation in content production, and (2) exploring the normative and operational implications of deploying such systems in live editorial contexts.

The key objectives are:

- To analyze the current technological stack used in AI-powered journalism, including NLP, NLG, and machine learning classifiers.
- To examine how news organizations are restructuring their editorial workflows to accommodate AI tools.
- To evaluate the ethical and legal concerns related to transparency, trust, and human oversight in automated journalism.
- To design and simulate a typical journalist's daily journey using an AI-integrated dashboard and content pipeline.
- To propose a conceptual architecture that balances automation with editorial control, considering current European regulatory frameworks.

1.3 Research Aim and Objectives

The primary aim of this research is to explore how AI and generative AI systems are transforming journalistic workflows, from data ingestion to content distribution. Specifically, this study examines both the technical architecture and the editorial consequences of integrating AI into newsrooms.

The objectives are as follows:

- To analyze the current technological stack used in AI-powered journalism.
- To investigate how news agencies are redesigning workflows to accommodate AI.
- To evaluate the ethical, editorial, and labor-related implications of automated journalism.
- To design and simulate a journalists journey within an AI-assisted newsroom.

1.4 Research Questions

Building on the objectives outlined above, this thesis addresses a series of interrelated research questions aimed at bridging both conceptual and applied gaps in the study of AI in journalism. While prior work has examined specific technologies or ethical implications in isolation, this research seeks to connect technical implementation, workflow design, and editorial values in an integrated model. The guiding questions are:

1. How is artificial intelligence currently implemented in news production workflows?
2. What roles do technologies such as NLP, NLG, and machine learning play in these implementations?
3. What are the key benefits and risks associated with AI in journalism, particularly in relation to editorial control, transparency, and efficiency?
4. How can generative AI enhance or compromise journalistic standards and integrity?
5. What design principles can guide the responsible integration of AI into newsroom infrastructures under real-world constraints?

These questions aim to ground the discussion in both theory and practice, ultimately contributing a systems-level perspective on humanAI collaboration in editorial contexts.

1.5 Methodological Overview

Given the multidisciplinary nature of the topic, this study adopts a mixed-methods approach that integrates technical inquiry with critical analysis. The goal is not only to evaluate what is technologically possible but also to understand how those technologies reshape the social, ethical, and operational dimensions of journalism.

The methodology consists of the following components:

- **Literature Review:** A systematic synthesis of academic and industry sources to trace the development of AI applications in journalism and identify research gaps.
- **Case Studies:** Analysis of real-world implementations (e.g., AP, BBC, Schibsted) to understand how AI systems are deployed, adapted, and overseen in operational newsrooms.
- **System Simulation:** A conceptual pipeline modeling an AI-powered newsroom workflow, including stages such as data ingestion, NER, summarization, and generative rewriting.

- **Critical Analysis:** An evaluation of the societal and normative implications of AI integration in journalism, with a focus on transparency, labor dynamics, and editorial accountability.

This combination allows for a holistic exploration of AI in journalism, balancing empirical observations with design-oriented thinking.

In the development of this thesis, I made use of generative AI tools, to support drafting, text refinement, and language editing in English. These tools were also used in the early stages to explore phrasing alternatives and synthesize technical content. All final material has been critically reviewed and validated by the author.

1.6 Structure of the Thesis

The thesis is organized to progressively build the case for responsible AI adoption in newsrooms, moving from context and problem definition to applied solutions and ethical synthesis:

- **Chapter 1 Introduction:** Introduces the context, defines the research problem, outlines objectives, and frames the methodology.
- **Chapter 2 Literature Review:** Synthesizes the state of the art in automated journalism, covering both technical foundations and critical debates.
- **Chapter 3 The Problem:** Analyzes structural, editorial, and economic pressures that are driving news organizations toward automation.
- **Chapter 4 The Solution:** Describes AI-driven tools and architectures currently used in journalism, with a focus on workflow automation and generative models.
- **Chapter 5 Designing the Journalist Journey:** Proposes and simulates a realistic day-in-the-life of a journalist using an AI-enhanced editorial dashboard, offering a concrete vision of humanAI collaboration.
- **Chapter 6 Discussions and Conclusions:** Revisits research questions, synthesizes findings, explores ethical and legal implications, and reflects on the broader contributions of the thesis.

This structure ensures a coherent narrative arc, moving from theoretical framing to applied modeling, and culminating in critical reflection and future outlook.

Chapter 2

Journalism in the Digital Era

2.1 Defining Automated Journalism

2.1.1 Conceptual Foundations

Automated journalism also known as algorithmic or robot journalism refers to news stories created by computer software with little to no hands-on involvement from humans after the initial setup. At its core, this approach relies on pre-written templates and structured data feeds to generate straightforward, factual reports (Danzon-Chambaud, 2021; Graefe, 2016). As Graefe describes, once the algorithm is in place, it can take over the entire production chain: gathering data, analyzing it, writing the story, and even publishing it all without further human input (Graefe, 2016). The key point is that these are fact-based, data-driven stories, produced at scale using automation. Algorithms sift through structured datasets to pull out meaningful bits like names, stats, or events and slot them into pre-set narrative formats or use language models to craft the story. You'll see terms like *algorithmic journalism* and the more ambiguous *robot journalism* pop up in the literature (Danzon-Chambaud, 2021), but the common thread is the use of natural language generation (NLG) to produce readable, coherent news content. Essentially, this is computational journalism where much of the traditional reporting process gets handed over to software, which uses natural language tools to turn raw data into finished articles.

Researchers often point out that automated journalism is built on the broader field of computational journalism (Dalgali and Crowston, 2020; Diab, 2023). It reframes certain routine news tasks like sports recaps, financial updates, or weather reports as problems that can be solved with code. Because it's heavily data-dependent, this method works best on topics with lots of structured info like company reports, game stats, election tallies, and so on. In that sense, it's a continuation of traditional data journalism, just with machines now taking over the repetitive grunt work that junior staff might have handled in the past. And by its very nature, once the system is live, it doesn't require much (or any) manual writing (Graefe, 2016). As one industry overview puts it, these systems can automatically generate news stories without human intervention after the initial programming of the algorithm (Press, 2024). Despite being machine-produced, the output is designed to read like standard news since fluency and style are either pre-defined or learned. However, the responsibility for content quality still rests with people. Editors and developers are the ones who set the templates, choose the data, and take the blame when something goes wrong (Danzon-Chambaud, 2021; Graefe, 2016). So at the end of the day, automated journalism is

Latest From This Author



CALIFORNIA

3.9 earthquake centered rattles San Francisco Bay Area

Moderate shaking was reported with a magnitude 3.9 earthquake near Dublin, Calif. Residents in San Francisco, Fremont and Richmond reported weak shaking.

March 17, 2025

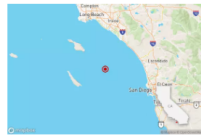


CALIFORNIA

Magnitude 2.9 earthquake registered in Los Angeles

A small temblor shook at 3:59 p.m. Thursday less than a mile from South Pasadena, according to the U.S. Geological Survey.

Oct. 31, 2024



CALIFORNIA

Series of small earthquakes rumbles off SoCal coast

A magnitude 3.6 earthquake struck in the waters off Catalina Island on Sunday evening on the heels of two smaller quakes in the same area.

Oct. 20, 2024



CALIFORNIA

L.A. rattled by three more small earthquakes north of Malibu

All three earthquakes were reported within the same general area north of Malibu where a magnitude 4.7 temblor had hit four days earlier.

Sept. 16, 2024

Figure 2.1: Today's Quakebot online published content - LA Times

really about programming the news, while leaving humans in charge of editorial standards and accuracy (Fearn, 2025; Hasan et al., 2023).

2.1.2 Historical Evolution and Trends

The idea of computers writing news is not brand new, but practical adoption accelerated in the 2010s. One of the earliest high-profile cases was the Los Angeles Times *Quakebot* in 2014, which automatically published earthquake reports within minutes of seismic events (Barca, 2022). Around the same period, major news organizations began partnering with data firms and AI companies to automate routine beats. For example, the Associated Press (AP) rolled out an automation program in 2014 to cover corporate earnings: using data from research firms and software from Automated Insights, AP began generating thousands of earnings summaries each quarter (Graefe, 2016; Press, 2016). This yielded a dramatic increase in coverage; by some accounts AP was producing roughly 3,700 company-earnings stories per quarter by 2015, a volume far beyond what its human reporters alone could sustain (Graefe, 2016). Encouraged by this success, AP and others expanded into sports: in 2016 AP announced it would automate coverage of all minor-league baseball games (142 teams across 13 leagues) using Automated Insights Wordsmith platform (Press, 2016) (Previously AP had covered only a handful of games manually.) The same technological wave hit other countries: for instance, Le Monde in France launched an automated system for election results in 2015 (with local AI startup Syllabs), and Swiss publisher Tamedia used algorithmic scripts for referendum and election coverage in 2019 (Danzon-Chambaud, 2021; Fanta, 2017). The BBC and Reuters have also experimented with similar systems (e.g., for U.K. election night coverage).

Academic reviews note that research interest in automated journalism also took off post-2010. Early scholarly articles (K. Dörr, 2016; van Dalen and, 2012) predicted major shifts, and by the mid-2010s

commentators were actively debating its implications. The World Editors Forum even named automated news-writing a top trend for newsrooms in 2015 (Graefe, 2016). In the research community, the field overlapped with studies in computational journalism and media automation. For example, Graefe (Graefe, 2016) and Diakopoulos (2019) examined how pre-set templates and machine learning could handle repetitive news tasks, while Gynnild (2014) and others explored new skills needed by journalists. Importantly, literature surveys emphasize that most early scholarship has centered on large, data-rich media organizations in Europe and North America (Danzon-Chambaud, 2021). There is less evidence on how non-Western newsrooms or smaller outlets use these tools. Nonetheless, a clear pattern has emerged: automated journalism grew from niche projects (like sports recaps and earthquake alerts) into a recognized mode of production, with dozens of major outlets now incorporating some form of AI-driven reporting (Danzon-Chambaud, 2021; Fanta, 2017). In summary, the historical trajectory shows a shift from experimentation to integration: automation began with isolated cases in the early 2000s, became commercially viable in the 2010s, and by the early 2020s is a routine element in many news organizations toolkits.

2.2 Technological Foundations

Automated journalism is underpinned by several key AI technologies. These range from natural language processing and generation to traditional machine learning and data engineering. We consider each class of techniques and how it contributes to the automated-news pipeline.

2.2.1 Natural Language Processing (NLP): Extraction and Summarization

Natural Language Processing (NLP) plays a key role in turning raw data and text into usable material for automated news writing. In this context, NLP tasks like named-entity recognition, event detection, and sentiment analysis help break down massive datasets like company earnings reports or sports stats into clean, structured inputs that a generation system can use. Take a press release or a financial filing, for example: an automated system can scan the document, pull out important figures or developments, and label them accordingly. More advanced techniques even allow for summarization, letting algorithms distill long documents into shorter abstracts or highlight key takeaways. While this kind of summarization isn't always essential for rigid, template-based reporting, it becomes crucial in workflows where humans are still involved or where the goal is to turn technical data into plain, readable summaries. Essentially, NLP acts as a bridge that turns messy, unstructured input into the building blocks a news-writing algorithm can work with.

That said, most real-world automated journalism tools today use fairly basic NLP. Some systems are little more than scripts that extract numbers from a database and slot them into blanks in a pre-written story template (Graefe, 2016). For instance, a bot covering sports might grab the final score and highlight a few top players, then insert that info into a fixed game summary format. On the more advanced side, some newsrooms are exploring summarization tools that can turn press releases into brief bullet points, or scan social media for signs of breaking news. The field of summarization is complex with methods like extractive vs. abstractive models but for automated journalism, what really counts is accuracy. Summaries need to be factually correct and maintain a professional, journalistic tone. Current research shows that even powerful neural summarizers often need fine-tuning or editorial

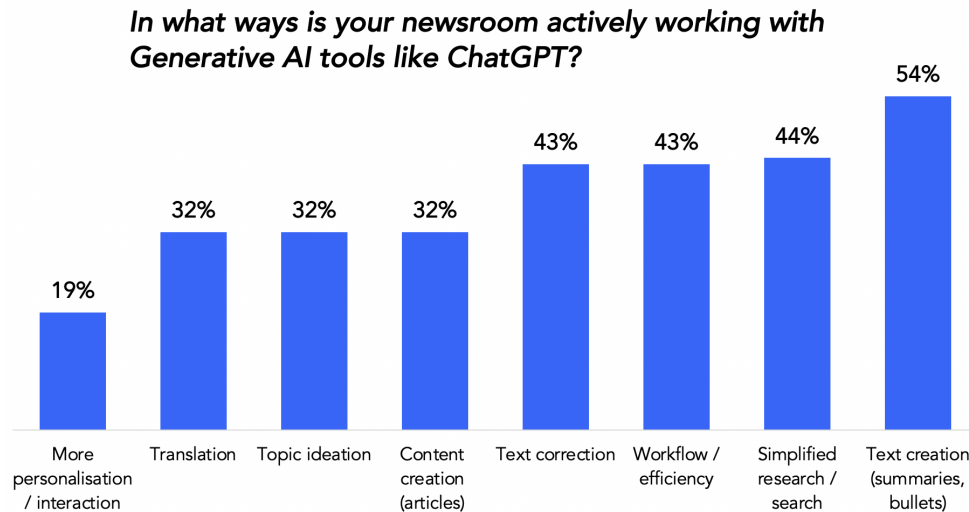


Figure 2.2: NewsRooms usage of AI Tools
Source: WAN-IFRA, 2023

oversight to meet newsroom standards (Narayanan, 2018). So in most workflows, NLP is used up front to filter and interpret the data making sure that what gets passed on to the generation step is relevant, accurate, and ready to be shaped into a story.

2.2.2 Natural Language Generation (NLG): Report Writing and Rewording

Natural Language Generation (NLG) sits at the core of automated journalism. It's the part that turns structured data into readable stories. These techniques range from simple systems that fill in blanks in a template, to more sophisticated neural models that generate fluid, varied prose. Most of today's automated news tools still stick pretty close to the template-based end of that spectrum. In these setups, editors or data scientists design reusable story outlines with placeholders for key facts. The system then plugs in the data and applies some basic grammar rules to generate a clean, consistent output. For instance, an NLG engine might use a template like [TeamA] defeated [TeamB] with a score of [X][Y], pulling the values straight from a sports database. This method helps keep the output grammatically sound and factually reliable, but it often results in fairly repetitive, formulaic writing.

That said, there's growing interest in going beyond these rigid templates. New research and machine learning techniques are helping algorithms learn how to shape stories more flexibly. One study, for example, describes a data-driven architecture that can generate election coverage in multiple languages (Leppänen et al., 2017). These more advanced systems might use statistical models or even large neural language models to paraphrase and rephrase content, making the final stories feel a bit more natural. But there's always a trade-off: journalistic content has to prioritize accuracy and transparency. That's why even the more creative NLG tools in newsrooms are built with safeguards. Often, they generate a draft that still gets reviewed and polished by a human editor. In short, automated journalism uses NLG to transform raw data into full-length stories, with the sophistication of the language generation tuned to fit the editorial goals and trust standards of the organization (Graefe, 2016).

2.2.3 Machine Learning Techniques: Classification, Clustering, Predictive Models

Beyond NLP and NLG, machine learning (ML) also plays a supporting role in automated journalism. Classification algorithms, for instance, can help sort incoming information like organizing news tips by topic or automatically tagging article themes. Clustering techniques might group related stories or data points, which is useful for spotting trends or avoiding duplicate coverage. Predictive models, on the other hand, can be used to flag likely future events like forecasting election outcomes or sales numbers based on previous patterns.

In real-world applications, most of these ML tasks rely on training models with historical news data or external datasets. For example, a classifier could be trained to spot press releases about corporate earnings and flag them for automated reporting. It could also learn to tell the difference between routine and unusual events. Similarly, clustering methods might alert a system that multiple traffic accidents have occurred in the same area possibly pointing to a broader story. While published systems tend to highlight template-based NLG more than ML components, the academic research shows that machine learning is doing a lot of the behind-the-scenes work. According to field reviews, ML helps journalists manage the overwhelming volume of data, freeing them up to focus on more complex or sensitive angles that require human judgment (Carlson, 2015; Dalgali and Crowston, 2020).

2.2.4 Web Scraping, Data Ingestion, and Real-Time Feeds

All of the technologies mentioned so far depend on one crucial ingredient: data. In practice, automated journalism runs on a strong, well-built data pipeline. Newsrooms tap into a wide range of data sources everything from structured databases (like sports stats, financial figures, weather sensors, or public records) to real-time feeds such as stock tickers, election results, and even social media APIs. Some systems also rely on scraped content from websites or social platforms.

Many automation workflows start with APIs. For example, the Associated Press's baseball bot pulls real-time game statistics straight from Major League Baseball's data services (Press, 2016). In other situations, developers or journalists write custom web scrapers that grab info from online sources at regular intervals. Think of a bot that refreshes an official election results page every few minutes during vote counting. Real-time data triggers can also come into play like the LA Times earthquake bot, which watched the USGS earthquake API and auto-generated a story whenever a quake above a certain magnitude hit.

At the end of the day, data ingestion is the first and most essential step in the pipeline. For automated journalism to work, there needs to be a steady stream of clean, structured data flowing into the system. Once the data is in, it's handed off to NLP tools for analysis, and then to NLG models to generate the final news copy. As Graefe (2016) puts it, what sets automated journalism apart is the ability to automate each step of the news production process from gathering data all the way through to publication (Graefe, 2016).

2.3 Benefits of AI in Newsroom Workflows

Integrating AI into news production offers several potential advantages. These are often highlighted by both industry advocates and scholars studying automated journalism.

2.3.1 Efficiency and Speed

The most obvious advantage of automated journalism is speed. Once triggered, these systems can generate and publish news stories almost instantly. For breaking events driven by data, AI can pull together a complete report far faster than any human. More broadly, research has shown that automated news tools can provide all the facts a human reporter would, but much faster and at a lower cost (Graefe, 2016).

This efficiency shows up clearly in newsroom output. After the Associated Press adopted Automated Insights platform in 2014, it began producing over 3,500 earnings stories every quarter about ten times the volume it had managed manually. In sports, too, AI bots are now generating hundreds of game summaries each week, ensuring that local teams get timely recaps and fans stay instantly informed. This kind of productivity boost frees up human journalists for more complex work. Once routine coverage is automated, newsrooms can shift their focus to investigative reporting, feature writing, or editorial projects.

Those managing the automation systems take on new responsibilities as well. They act as overseers reviewing output, checking for errors, and refining the templates and rules the AI follows. But the bigger picture is clear: automation helps newsrooms do more with less time (Carlson, 2015). By offloading formulaic, data-heavy reporting to algorithms, news organizations can provide round-the-clock coverage and frequent updates without needing to drastically expand their teams.

2.3.2 Scalability and Personalization

Closely tied to speed is the advantage of scale. AI systems can handle far more stories and data points than any human team could manage. While human reporters might only cover high-profile games or major companies, an automated system can crank out thousands of reports many of them hyper-local or niche. A good example is the APs baseball bot, which now produces summaries for every minor league game, giving visibility to dozens of teams that previously went uncovered. Tools like Wordsmith make this even easier: a single algorithm can operate across different domains finance, sports, weather just by hooking into new data sources. In theory, that means a newsroom could massively expand its coverage without needing to hire a matching number of new reporters.

Another major benefit is personalization. Since AI-generated content is built from structured data and rules, it can be easily tailored for different audiences or formats. Algorithms can emphasize different angles depending on who the reader is or where the content appears. For instance, a story about a companys earnings might focus on profit margins for investors but highlight job impacts for a general audience. In fact, some organizations already use AI to customize newsletters or alerts, sending slightly different summaries to each reader based on preferences or Browse history. Graefe (2016) highlights

this potential, noting that automated journalism can personalize [content] to the needs of an individual reader without compromising on factual accuracy (Graefe, 2016). This ability to personalize stems from the software's inherent flexibility: once the basic structure of a story is in place, variations can be generated at scale. So, in short, AI allows journalism to be both broader in reach and more targeted in tone—something traditional methods struggle to achieve efficiently or affordably.

2.3.3 Innovative Content Forms

Beyond efficiency, AI has also opened the door to entirely new types of content formats that would've been impractical or impossible before. Automated tools are now being used for more interactive and experimental forms of storytelling. One standout example is the rise of news quizzes and chatbots. In 2023, BuzzFeed made headlines by announcing it was using AI to create personalized quizzes and even some travel guides, aiming to boost engagement with content that feels more like a game than traditional news. Similarly, several outlets have tested chat-style bots that users can ask for real-time updates like the latest sports scores or financial stats.

Another emerging area is synthetic media. In a notable example, Chinese state media developed AI-powered virtual anchors—digital presenters who can deliver the news 24/7. These systems are still fairly new and typically used in specific scenarios, but they point to a broader trend: AI isn't just about writing; it's about expanding how stories are told. Algorithms are now being used to automatically generate visuals like charts, infographics, and even audio narration. In fact, several AI services can now turn a written news article into a podcast on demand.

At the core of all this is the idea of a more dynamic news experience. Imagine a live dashboard that talks you through the numbers as they update, or an augmented-reality bulletin that shifts as the story evolves. While the academic literature on automated journalism is still just scratching the surface of these possibilities, the early examples are promising. What's clear is that AI adds creative range to journalism—from fun, personalized quizzes to immersive, multimedia storytelling—by making it easier to produce this kind of content at scale.

2.4 Risks and Challenges

Alongside benefits, automated journalism raises significant concerns. Scholars and journalists have identified ethical and practical risks that must be managed.

2.4.1 Data Bias and Algorithmic Fairness

One of the biggest concerns around AI in journalism is the potential for bias. These systems can easily inherit and even amplify the prejudices baked into their training data. If the datasets or rules behind automated news algorithms reflect past inequalities or narrow viewpoints, the content they generate may unintentionally sideline certain groups or reinforce stereotypes. For instance, if a sports bot's language templates are developed by a culturally uniform team, they might fail to account for the diversity of athlete names or use phrasing that comes off as insensitive. As (Fearn, 2025) points out, underlying

datasets are often biased because the people entering the data bring their own perspectives and these can be skewed by demographic imbalances or unconscious bias.

In the context of journalism, this presents a real risk: automated tools could unintentionally ignore minority communities or repeat problematic narratives. Fearn warns that without active oversight, AI systems might reinforce existing inequalities in the way news is reported (Fearn, 2025). Bias also shows up in subtler ways like which stories get automated in the first place. Algorithms typically gravitate toward topics that are rich in structured data, such as finance or sports. That means more complex social issues where data is messy, incomplete, or non-existent might get sidelined. The result is an editorial skew toward beats that are easy to quantify, potentially leaving out stories that require more context or human nuance.

Researchers note that even small algorithmic choices, like which keywords to monitor, carry implicit value judgments. So ensuring fairness isn't just a technical task; it's also an editorial one. That's why many in the field stress the need for regular audits of both datasets and output. Newsrooms need to ask: which voices are missing? Are we reinforcing clichés? What assumptions are built into the system? Human oversight is key; editors must remain involved to catch and correct biases before they reach the audience (Fearn, 2025).

2.4.2 Transparency and Editorial Accountability

Closely tied to the issue of bias is the question of transparency and accountability in automated journalism. Because algorithms can produce and publish stories without direct human writing, it becomes harder for both audiences and even newsroom staff to understand how a piece was created. In traditional reporting, the chain of responsibility is clear: a reporter writes the story, and an editor reviews it. But when an algorithm is the "author," there's no obvious person to hold accountable for mistakes or misjudgments.

That's why industry experts stress that responsibility for automated content must still fall on humans. As Graefe explains, algorithms themselves can't be held legally or ethically accountable, so it's the individuals who design, deploy, or oversee these systems—data journalists, editors, developers—who must take ownership of the outcomes (Graefe, 2016). This raises the need for stronger transparency. News organizations are increasingly called to adopt and publicize clear policies, often called algorithmic principles or usage guidelines that tell audiences when a story was machine-generated and how the system behind it works.

Some media ethicists argue this kind of transparency should be treated like traditional source attribution: readers should know what parts of an article were written by a machine and which datasets were used in the process. Without that clarity, there's a risk of undermining trust, especially if readers start to suspect that stories are being produced behind the scenes by black-box systems. And as personalization becomes more common, these concerns grow. Tailored news feeds, driven by algorithmic curation, can create so-called filter bubbles (concept that we will cover later) if users don't realize their content is being selectively filtered (Pariser, 2011).

In short, automated journalism isn't just about speed and efficiency; it also requires new ethical norms.

Transparency and explainability need to be built into the process (Carlson, 2015; Graefe, 2016). Newsrooms must make sure that both humans and algorithms follow shared editorial standards, and that any errors in automated stories can be quickly identified, traced, and corrected by a real person.

2.4.3 Impact on Labor and Editorial Roles

One of the ongoing concerns around automated journalism is its potential impact on employment in the industry. Many people worry that bots will eventually replace human reporters, leading to widespread job losses. But in practice, the picture has been more nuanced. So far, surveys and case studies suggest that automation hasn't triggered mass layoffs in newsrooms. In fact, Dalgali and Crowston cite AP executives who report that their algorithmic journalists haven't displaced any human staff (Dalgali and Crowston, 2020). Instead, many organizations are reallocating staff to more strategic or creative roles. As one insider explained, the time saved through automation often gets reinvested either in expanding coverage areas or in producing deeper, more analytical content (Press, 2016).

Some publishers have even created new positions to manage automation, such as automation editors who oversee bots, tweak templates, and maintain quality control. For journalists with skills in data analysis or coding, this shift can open up exciting new opportunities. These roles act as a bridge between the editorial team and the algorithms powering automated content.

That said, integrating AI into the newsroom isn't without challenges. It can create stress and uncertainty among staff. Recent studies show that journalists who view AI as a threat to their jobs often experience heightened anxiety or even depression (Upadhyay et al., 2024). The concern isn't just about being replaced; it's also about being retrained. As newsrooms adopt automation, journalists may need to pivot into roles that involve supervising algorithmic output, managing data, or shaping templates instead of writing traditional stories.

There's also the issue of deskilling. If entry-level reporting tasks are fully automated, newer journalists might miss out on foundational experiences like basic reporting, fact-checking, or story development that help them grow in the profession (Posetti, 2018). Some scholars warn that, without careful planning, automation could hollow out the early-career pipeline and widen the skills gap between editorial and technical roles.

Still, there's a flip side. By handling repetitive work, AI could help make journalism more financially viable during tough times, possibly even preserving jobs that might otherwise be cut (Carlson, 2015). Much of the impact depends on how news organizations choose to implement these tools. Outlets that treat AI as an assistant rather than a replacement often see gains in both efficiency and output. But those using it purely to slash costs may end up reducing staff.

In short, automation is reshaping newsroom roles rather than eliminating them. It's creating new technical positions and shifting the focus of existing ones. With thoughtful integration, training, and support, AI has the potential to complement human reporters, not replace them.

2.4.4 Content Farms

Content farms are online platforms that produce large volumes of articles, often using Natural Language Processing (NLP) models to generate news-like content aimed at attracting web traffic and ad

Is there resistance to use Generative AI tools?

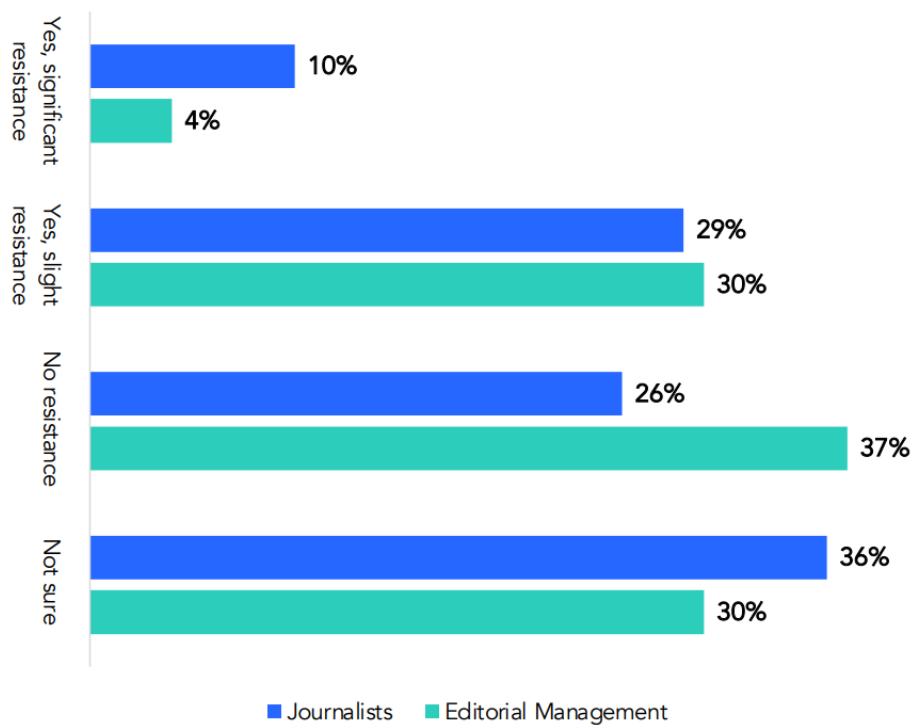


Figure 2.3: Stakeholders and GenAI
Source: WAN-IFRA, 2023

revenue. These websites are typically filled with synthetically generated AI texts. While they may not necessarily be part of coordinated misinformation campaigns, their sole objective is to maximize user visits(Puccetti et al., 2024).

Although this is a relatively recent phenomenon, the increasing capabilities of Generative AI are making content farms more disruptive. These sites generate revenue through traffic-driven advertising, capitalizing on sensational headlines and trending topics. Their business model relies on programmatic advertising, where ads are delivered by the ad-tech industry without considering the quality or trustworthiness of the websiteNewsGuard, n.d. As a result, even top brands may inadvertently fund these websites.

Content farms often adopt generic names and layouts to appear as trustworthy news sources and appeal to a broad audience. NewsGuard has been actively tracking such platforms and, to date, has identified 1,254 "Unreliable AI-Generated News" websites across 16 languagesNewsGuard, n.d., as shown in Figure 2.4. These websites meet the following four criteria:

1. A substantial portion of the content is generated by AI.
2. There is clear evidence of minimal or no human oversight (e.g., error messages or chatbot-specific phrasing).
3. The site is presented using a generic layout and a benign or generic name.

4. The site does not clearly disclose that its content is AI-generated.

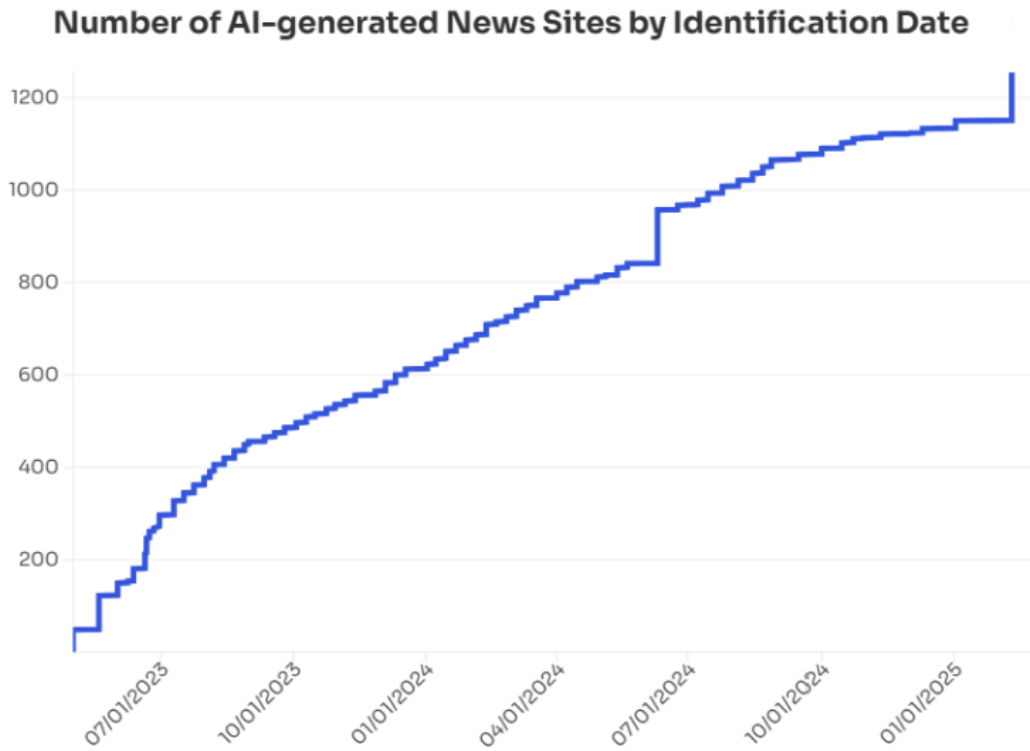


Figure 2.4: Number of AI generated News Sites by NewsGuard

Detecting whether a text is AI-generated remains a significant challenge for users, while at the same time, building effective content farm models has become relatively simple. This issue is illustrated in an insightful study by Puccetti et al. (Puccetti et al., 2024), where the researchers fine-tuned a relatively old LLM, LLaMA-65B, using a dataset of 40,000 Italian news articles. The resulting model was capable of misleading native Italian speakers, who were only able to distinguish between human- and AI-generated texts with an accuracy of 64%.

The study also shows that existing AI-detection tools, which rely on token likelihood estimation or supervised classification, outperform human judgment. However, these tools present several practical limitations. They often require access to internal model probabilities or large labelled datasets, making them resource-intensive and not easily deployable in real-world scenarios.

2.5 Summary of Literature Gaps

Despite the surge of interest in automated journalism, the academic literature still leaves plenty of questions unanswered. Reviews of the field highlight that much of the research so far has been exploratory or descriptive, with relatively few in-depth, empirical studies on long-term impacts (Danzon-Chambaud, 2021). One major gap is the lack of international perspective. Most case studies center on U.S. and European newsrooms, leaving us with limited insight into how automation is unfolding in regions like Asia, Africa, or Latin America.

There's also a call for stronger theoretical frameworks. As Danzon-Chambaud (2021) argues, we need to analyze automated journalism through institutional and organizational lenses drawing on concepts like field theory or Bourdieu's work to better understand how these technologies reshape power dynam-

ics within newsrooms (Danzon-Chambaud, 2021). Another underexplored area is evaluation. While there are anecdotal reports about increased productivity, we still lack systematic comparisons between machine-generated and human-written news, especially regarding quality, tone, or long-term audience trust.

Similarly, very few studies dig into the broader societal implications. Issues like misinformation, public discourse, and reader engagement are often mentioned, but rarely investigated in depth. Ethical discussions are also still in their early stages. While there's growing interest in building norms around transparency and fairness in AI-generated journalism, comprehensive guidelines remain scarce.

In short, the literature has done a solid job of identifying what automated journalism is its tools, terminology, and early use cases but it still has room to grow in terms of analytical depth and global scope. Future research could help close these gaps by focusing on longitudinal studies in real newsroom settings, exploring non-Western implementations, and embracing multidisciplinary approaches that merge technical insight with media theory.

Chapter 3

The Problem: Why Agencies Are Turning to AI

3.1 Changing Information Ecosystem

The modern media environment never sleeps. With the rise of the 24-hour news cycle and an overwhelming flood of content, audiences now expect constant updates and newsrooms are under pressure to deliver. As recent analysis puts it, we live in a fragmented media environment with seemingly endless sources of information (Center, 2024). To keep up, news organizations have had to significantly ramp up their output. A striking example is the already mentioned "The Associated Press": after introducing algorithms to generate earnings reports for small companies, it began producing around twelve times more stories than before far more than human reporters alone could handle. Numbers like that show just how much the 24/7 news cycle pushes outlets to scale fast. Any lag risks losing audience attention in a space crowded with competitors.

At the same time, digital platforms and mobile devices have exploded the number of sources vying for readers' eyes. Traditional news outlets are no longer just competing with each other; they're up against blogs, social media, wire services, and waves of user-generated content. Every day, readers are hit with dozens of headlines from across the internet, from legacy media to viral TikToks. To stand out in that noise, newsrooms often have to repackage and repost content across various platforms. They also optimize for search engines and social feeds, but that comes with a catch: they're increasingly dependent on platform algorithms that decide what content users see. So, no matter how reputable a news brand is, it still has to keep feeding the digital content machine or risk being drowned out altogether.

3.1.1 24/7 News Cycle and Content Saturation

Digital technology has completely reshaped the way news is published. In the pre-digital era, a newspaper or TV network could sit on a story until the next scheduled edition. Today, even a short delay risks losing the audience to a faster competitor. This has changed the game now, even minor updates or routine events are reported instantly, adding to the sheer volume of news that needs to be produced and managed. Journalists are expected to keep feeding the non-stop global news cycle, contributing to what many call content saturation. Readers are overwhelmed with information, and newsrooms are under pressure to constantly generate fresh stories just to keep up.

Take financial journalism, for instance. The Associated Press now uses an automated system that

WHY HAVE YOU STARTED ADOPTING AI TECHNOLOGIES?

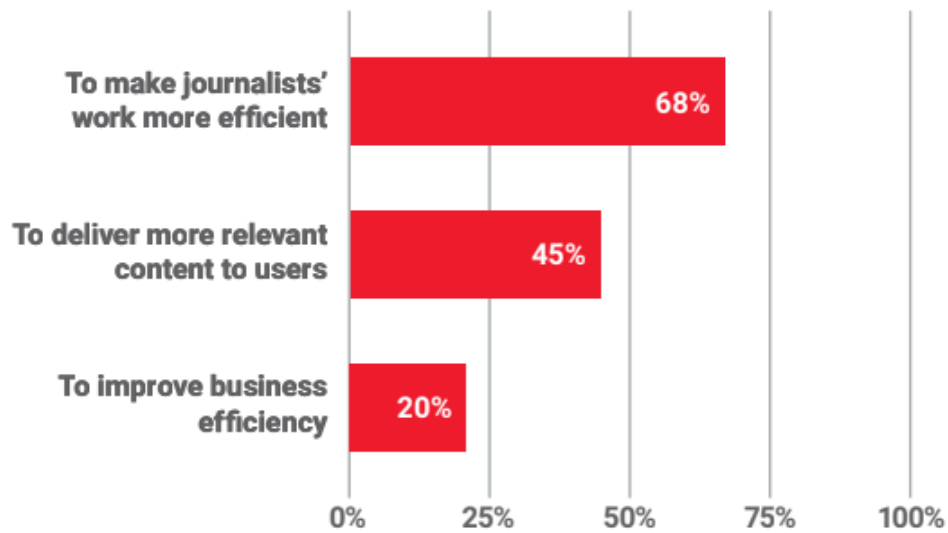


Figure 3.1: Reasons why adopting AI in newsrooms
Source: (Linden et al., 2021)

pulls data directly from company filings and quickly turns it into earnings stories. This approach has resulted in a twelvefold increase in coverage, especially of small companies that previously got little to no attention from human reporters. As the APs global business editor explained, automation has made it possible to cover far more firms than beforeexpanding reach without overwhelming staff. This highlights a bigger trend: as content demands soar, its increasingly unrealistic for newsrooms to rely solely on human writers. To keep pace, many are turning to algorithmic tools that can scale output and meet the growing appetite for news.

3.1.2 Rise of Data-Driven and Personalized Journalism

At the same time, journalism has grown increasingly data-driven and personalized. With structured data now widely availablethink financial stats, public records, sports scores, or polling resultsnewsrooms have gained powerful tools for both storytelling and automation. Its now common for data journalism projects to dig through large datasets, such as government archives or national surveys, to uncover trends or create interactive graphics. Alongside this, many news organizations are also using audience analytics and recommendation algorithms to tailor content to individual readers. For example, online platforms often track behaviorclicks, scroll depth, reading time, and sharesand use machine learning models to recommend articles or even customize homepage layouts for different users (Bradshaw and Rohumaa, 2011).

This growing use of data and personalization reflects how media outlets are adapting to evolving audience habits. With so many choices out there, keeping reader attention has become a challengeand algorithms help make content more engaging and relevant. But this approach also introduces new editorial concerns. The concept of the algorithmic filter bubble highlights how personalization can narrow a readers exposure to diverse viewpoints (Pariser, 2011). And this isn't just theoretical: over half of U.S. adults now get at least some of their news through social media (Center, 2024), where platform

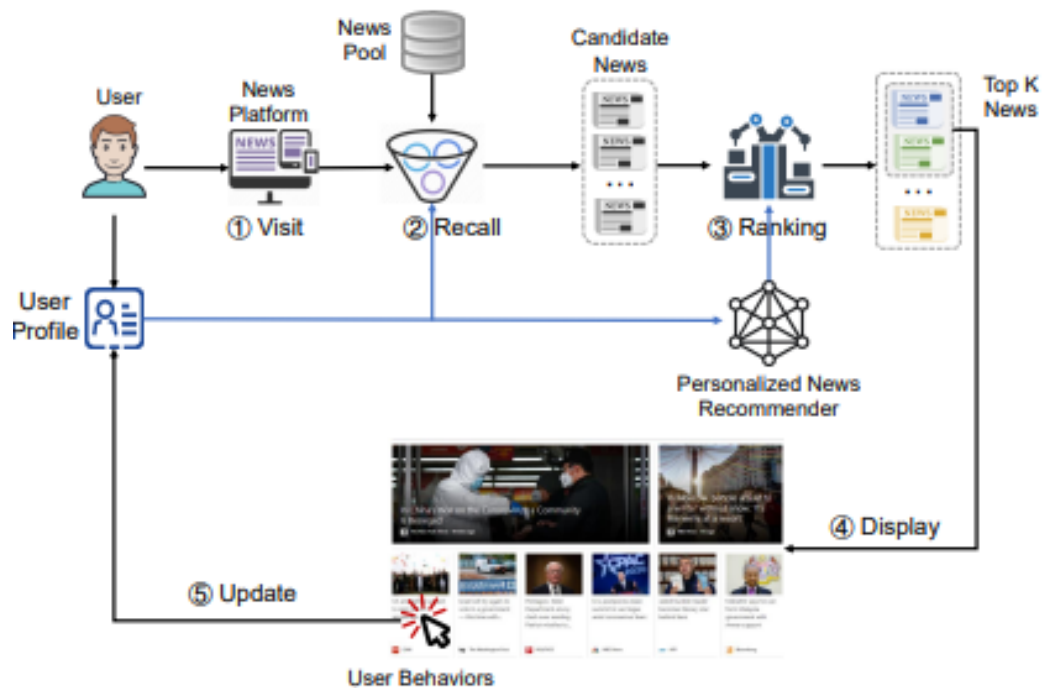


Figure 3.2: The workflow of news recommendation system
Source: Wu et al., 2021

algorithms decide what users see. That means readers may mostly encounter stories that confirm their interests or past behaviorslimiting the range of perspectives in their news diet.

In today’s fast-paced and information-rich environment, users are overwhelmed by the sheer volume of news published daily, making it difficult to identify content that best suits their interests. Recommendation systems help mitigate information overload and enhance the user experience (M. Li and Wang, 2019). Figure 3.2 provides a simple overview of how these systems work. When a user accesses a website, a set of news articles is generated as potential candidates for display. The personalized system ranks these articles based on the user’s interests, inferred from their previous interactions. The top-ranked articles are then shown, and the system tracks subsequent interactions, updating the user’s profile accordingly (Wu et al., 2021). In the absence of explicit user feedback (e.g., ratings), implicit actions, such as clicks or time spent on a particular article, are used to infer the user’s preferences.

Personalization is also an explicit expectation from users. According to a McKinsey report (Boudet and Vollhardt, 2023), 71% of users expect personalized content, and 76% express frustration when this expectation is not met. This insight is important as we later explore the potential advantages and drawbacks of this demandon one hand, enabling newsrooms to respond to market needs, and on the other, exacerbating concerns about "echo chambers."

So, while personalized content can boost engagement, it also raises questions about balance, diversity, and the role of journalism in an algorithmically curated world. For todays news organizations, this means walking a fine line: meeting audience expectations for relevance, while also protecting editorial integrity and social responsibility. In many ways, the rise of data-driven reporting and algorithmic distribution marks a fundamental shift in how journalism worksand how information reaches the public.

3.2 Bottlenecks in Traditional News Workflows

3.2.1 Manual Curation and Time Constraints

Traditional news production involves a lot of manual work at every step and under today's fast-paced conditions, that becomes a real bottleneck. Reporters and editors still need to sift through press releases, public records, social media posts, and wire service updates just to spot what's worth covering. Turning raw data or a quick event update into a finished article takes time: a journalist might need to verify stats, write the story, check for accuracy, and polish the final copy. All of that adds up, and much of it is repetitive. In reality, this means a significant portion of newsroom labor goes into routine information processing rather than high-impact journalism. As The Associated Press has noted, automating some of these tasks allows journalists to focus on more impactful aspects of their work. Without automation, reporters often find themselves bogged down by mundane chores like inputting data or formatting short updates, leaving less time for deeper reporting and analysis.

This challenge becomes even more intense when news is breaking. Editors need updates fast, often with verified quotes or confirmed details, but reporters don't always have the bandwidth to keep up. Most newsrooms don't have the staff to cover every story as it unfolds or to monitor every information stream in real time. That's why gaps appear. For instance, eyewitness reports might pop up on social media long before a journalist has time to verify and incorporate them into a piece. And since manually scanning these sources doesn't scale easily, real-time coverage can fall behind. In short, as the pace of news speeds up, the traditional workflow where people handle every step from scanning to writing to editing starts to slow everything down. For modern newsrooms, it's a structure that struggles to keep up with constant, global information flow.

3.2.2 Limitations in Real-Time Reporting

The demand for real-time updates adds yet another layer of pressure to traditional news workflows. Reporters and editors working against tight deadlines constantly have to juggle speed with accuracy. This often leads to a tough trade-off: publish fast with minimal verification, or take the time to fact-check and risk missing the moment. In practice, many newsrooms lean toward caution, which means slower story turnaround. But in today's digital environment, a delayed update can quickly get buried either by a faster competitor or by a viral post on social media. Audiences now expect news as it happens, and they're not inclined to wait. That leaves news organizations caught in a bind: move too quickly and risk errors, or move too slowly and lose relevance. And under traditional, manual workflows, there's rarely a perfect middle ground.

To keep pace, some outlets have turned to semi-automated tools. Algorithms can now scan social media in real time to spot trending topics or major developments and then alert journalists to investigate further. In some cases, broadcasters have integrated automated tickers or chatbots that push out live updates like election results or sports scores the moment new data becomes available. These tools don't replace human judgment, but they help speed up repetitive or time-sensitive parts of the reporting process. Still, many newsrooms continue to rely heavily on human monitoring and verification during breaking events, which stretches teams thin when several stories unfold at once.

Ultimately, the limits of manual reporting in a real-time world are becoming more apparent. To meet the growing demand for instant yet reliable updates, news organizations are increasingly looking to automation not to replace journalists, but to support them in delivering fast, accurate coverage under pressure.

3.3 Operational Pressures

3.3.1 Budget Constraints and Staff Reductions

Economic pressures have pushed many news organizations to cut costs, often through staffing reductions. Traditional sources of revenue like print and broadcast advertising have declined sharply, and the competition for digital subscriptions is fierce. As a result, many legacy media outlets have significantly downsized their newsrooms. In the U.S., for example, the Pew Research Center found that newsroom employment dropped by about 26% between 2008 and 2020 (Walker, 2021) as displayed in Figure 3.3. The losses were even more severe at newspapers, where editorial staff shrank by 57% during that time. Meanwhile, digital-native outlets have seen only modest growth, meaning the overall number of working journalists has declined substantially over the past decade.

With smaller teams and tighter budgets, newsrooms have fewer resources to cover the full range of stories. Routine updates or minor data releases may be skipped if journalists are stretched thin across core beats. In-depth investigations which take time and money often get sidelined. In this context, automation presents a practical solution. Once the infrastructure is set up, AI-generated content can be produced at a very low ongoing cost. Tasks like summarizing financial results, scanning public records, or generating brief news updates from structured data can be offloaded to algorithms. That frees up human reporters to focus on more complex or impactful stories.

In many ways, automation becomes a force multiplier; it helps newsrooms keep their output steady (or even grow it) despite having fewer staff. For organizations facing ongoing cuts, adopting AI tools isn't just an efficiency upgrade; it's a strategic necessity for maintaining basic coverage and staying afloat in a challenging media economy.

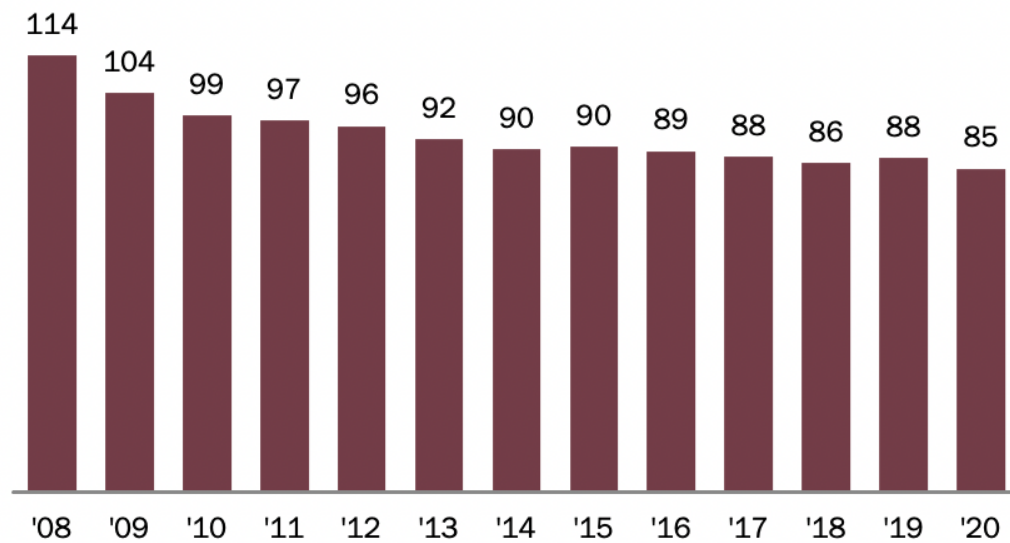
3.3.2 Competition with Non-traditional Publishers

Today's news consumers have more choices than ever, and many of them don't rely on traditional media outlets. Digital-native publishers, social media, blogs, podcasts, and even individual influencers now compete directly for attention. Platforms like Facebook and Twitter have become primary gateways to news: a 2024 survey showed that 58% of Americans prefer getting their news on digital devices like smartphones, tablets, or computers, compared to just 32% who still favor television and only 4% who turn to print (Center, 2024) as shown in Figure 3.5. Social media, in particular, now serves as a main news source for more than half of U.S. adults (Center, 2024). As a result, many people no longer visit traditional news sites directly; instead, they encounter stories through platform-curated feeds.

In this fragmented and fast-moving environment, legacy news outlets can't count on their reputation or established channels alone. To stay relevant, they need to publish content where their audiences already are: on social platforms, through mobile apps, and in formats optimized for on-the-go consumption.

Newsroom employment in the United States declined 26% between 2008 and 2020

Number of U.S. newsroom employees in news industries, in thousands



Note: The OEWS survey is designed to produce estimates by combining data collected over a three-year period. Newsroom employees include news analysts, reporters and journalists; editors; photographers; and television, video and film camera operators and editors. News industries include newspaper publishers; radio broadcasting; television broadcasting; cable and other subscription programming; and other information services, the best match for digital-native news publishers.

Source: Pew Research Center analysis of Bureau of Labor Statistics Occupational Employment and Wage Statistics data.

PEW RESEARCH CENTER

Figure 3.3: Newsrooms employment in USA - 2008-2020

Source: (Walker, 2021)

tion. Digital-native publishers often hold an edge here. They use powerful analytics, mobile-first formats like listicles and short videos, and viral marketing techniques to reach and engage younger audiences. Meanwhile, large tech platforms sort and serve news using algorithms, meaning traditional outlets must tailor their content to fit these systems whether that's by using trending keywords, structured metadata, or eye-catching headlines.

The competition for attention has pushed traditional newsrooms to speed up their production cycles and diversify how they deliver stories. In this context, automation offers a crucial advantage. Tools that can accelerate publishing, personalize content delivery, or adapt stories for multiple platforms help traditional outlets keep pace. To compete with newer, nimbler players, adopting automation isn't just helpful; it's becoming essential.



Figure 3.4: An example of Non-traditional competitor's content

A post from **Welcome to Favelas**, a popular Italian Instagram page that exemplifies how non-traditional publishers can reach large audiences with socially relevant content. Such platforms operate outside traditional media structures, often using humor, immediacy, and cultural resonance to inform and engage users posing a real competitive challenge to legacy news organizations. Their content is often reposted by traditional media.

3.4 Editorial Challenges

3.4.1 Quality Control at Scale

As news organizations begin automating parts of the editorial process, maintaining quality becomes a major challenge. Scaling up production with AI doesn't mean journalistic standards can be scaled down—accuracy, clarity, and fairness still have to be upheld. While automated tools are quite good at turning structured data into readable text, they often miss the nuance and context that a human reporter would naturally include. For example, an AI-generated article might correctly summarize a company's quarterly earnings, but skip over the significance of the numbers or fail to include important caveats. If these omissions aren't caught, they can lead to misunderstandings or even misinformation potentially damaging a newsroom's credibility.

In fact, a recent study suggests that audiences are more sceptical of news when they know it was produced with the help of AI Altay and Gilardi, 2024, as showed in Figure 3.6. Even when the facts are correct, the mere perception that a machine wrote the story can reduce trust. This shows that the challenge isn't just about getting the information right—it's also about maintaining public confidence in the process.

Another serious issue is the risk of algorithmic bias or error propagation. If the data feeding into an AI system is flawed or if the system's rules are poorly designed, those mistakes can ripple across everything it produces. For instance, if an algorithm leans heavily on historical data, it might unintentionally prioritize certain topics or perspectives, repeating past biases. That's why it's critical for editors to validate not just the outputs, but also the inputs and logic behind their automated systems.

Many newsrooms are addressing this by building in safeguards. These might include manual review of

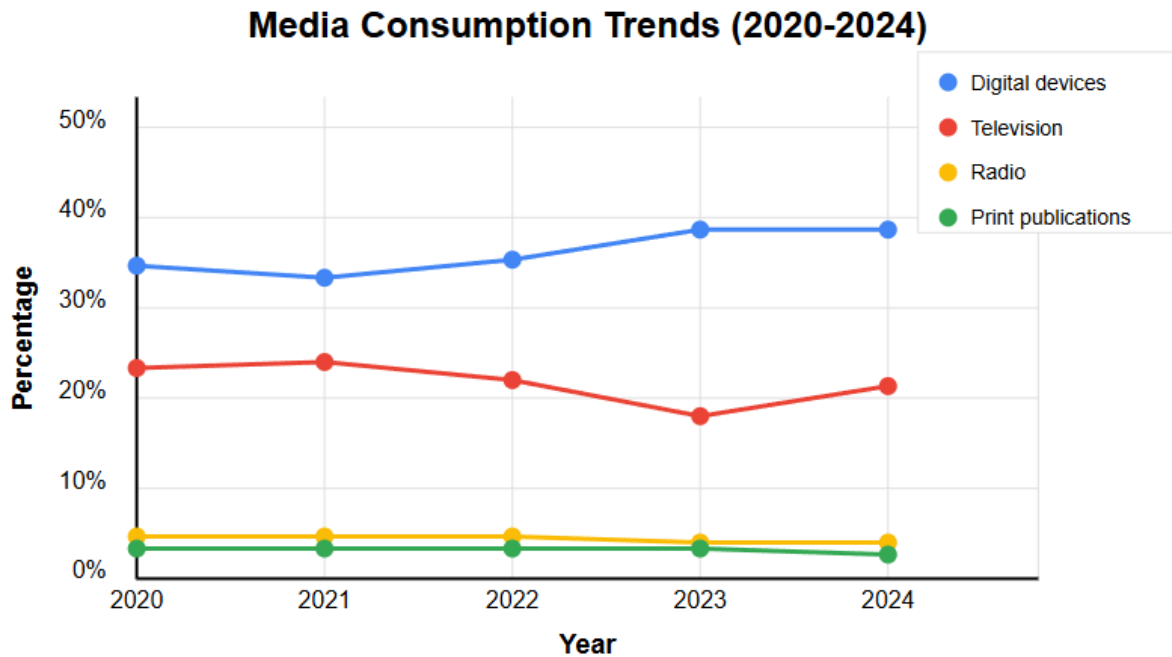


Figure 3.5: Media Consumption Trends
Source: (Center, 2024)

a sample of AI-generated articles, or setting up alerts for unusual or inconsistent results. As one expert puts it, automation can support journalism, but it requires careful journalistic governance (Johnson and Black, 2020). In practice, this means human editors need to be part of the process checking for logic gaps, spotting bias, and making judgment calls that AI simply can't. Maintaining high editorial quality at scale isn't just about having smart tools; it's also about building strong oversight workflows around them.

3.4.2 Maintaining Accuracy under Pressure

Accuracy is the foundation of journalism but under the pressures of speed and scale, it's becoming harder to safeguard. Automated systems, while efficient, can easily publish incorrect information if they're fed outdated or inaccurate data. For example, if an AI tool pulls numbers from an old database, it might present outdated figures as if they're current. And in the fast-paced world of real-time reporting, such mistakes can slip through unnoticed until after publication. Traditional fact-checking workflows often don't have time to keep up with the volume and speed of AI-generated content.

To tackle this, many newsrooms are adopting a hybrid approach. Automated drafts are typically reviewed by human editors before being published, and some organizations go a step further by deploying verification algorithms to cross-check facts. As (Diakopoulos, 2019) points out, human judgment is still critical. Most newsrooms maintain editorial oversight as a non-negotiable part of the process. The goal is to combine the speed of automation with the reliability of human review ensuring that content goes out quickly, but not carelessly.

In practice, this hybrid model involves new quality-control strategies. Some outlets run consistency checks or schedule random audits of AI-generated articles to catch errors that might otherwise go

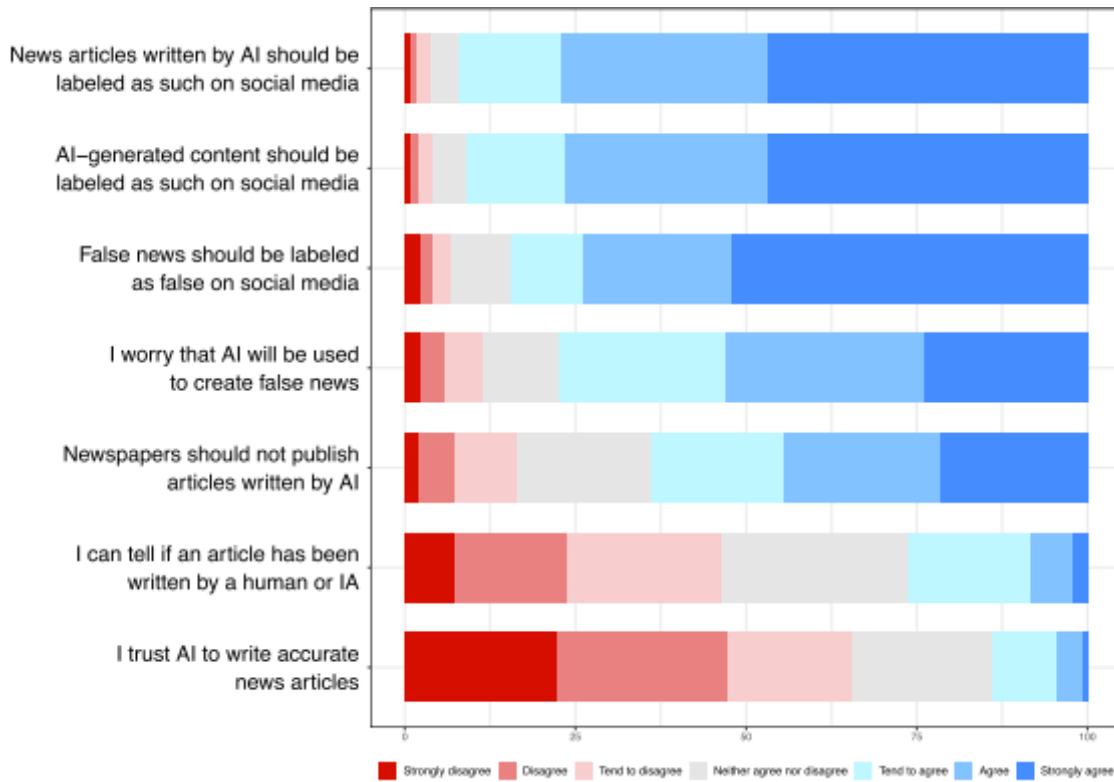


Figure 3.6: Attitude towards AI produced news
Source: Altay and Gilardi, 2024

unnoticed. While AI helps speed up production, it doesn't replace the editor's role; it reshapes it. The challenge is finding the right balance between efficiency and accuracy, so that automation supports journalism without undermining its core values.

3.5 Summary: Need for Workflow and Information Flow Automation

In summary, the modern news landscape presents multiple pressures that together make workflow automation a strategic necessity for many agencies:

- **Continuous content demand:** The 24/7 news cycle and vast digital environment mean audiences expect immediate, plentiful coverage. Traditional production processes alone cannot scale to meet this demand.
- **Resource bottlenecks:** Manual reporting and curation of routine information (e.g., data-driven news, breaking-event updates) create workflow delays. With limited staff and time, agencies cannot cover the volume of events and data inputs without technological assistance.
- **Economic constraints:** Shrinking budgets and staff cuts have forced news organizations to do more with less. Automation offers a way to maintain or expand output despite reduced human resources (Walker, 2021).
- **Competitive pressure:** Digital-native publishers, social platforms, and independent content creators challenge traditional outlets for audience attention. To keep audiences engaged, incumbents must adopt tools that enable rapid, targeted content production.
- **Editorial scaling:** Ensuring consistent quality and accuracy in a high-volume environment is dif-

DOES YOUR ORGANISATION HAVE A STRATEGY FOR AI?

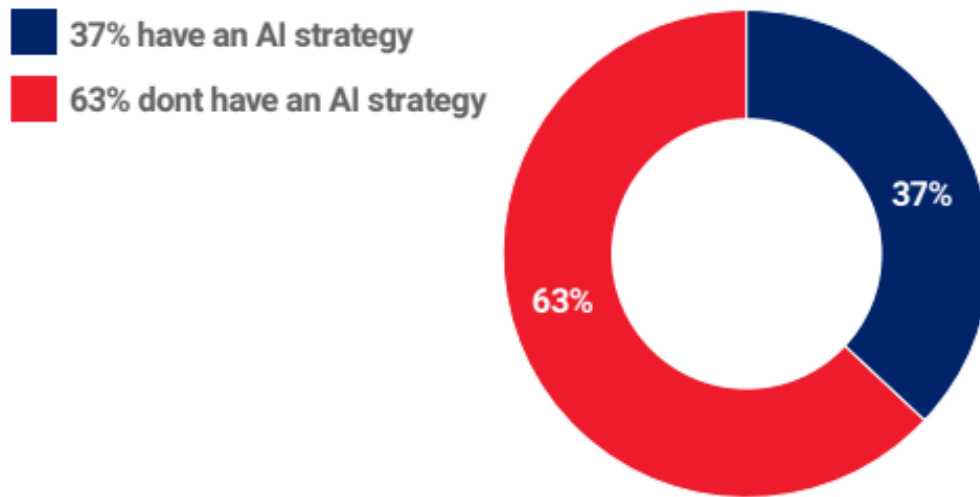


Figure 3.7: AI active strategies in newsrooms
Source: (Linden et al., 2021)

ficult. AI and automation present risks (e.g., errors, bias) that require new quality-control measures, but they also provide opportunities to strengthen workflows through algorithmic assistance (Appelman and Bien-Aimé, 2024).

Collectively, these factors create a clear incentive for news organizations to invest in workflow and information flow automation. The nature of todays information ecosystem — constant demand, instantaneous distribution, and data ubiquity — means that purely manual news production is no longer tenable. As (Diakopoulos, 2019) emphasizes, automation is becoming essential to sustain journalistic operations in the digital age. In the next chapters, we will explore how agencies are designing and deploying AI solutions to address these challenges.

Chapter 4

The Solution: Applying AI and GenAI to News Automation

4.1 Overview of AI Solutions in the News Industry

Media organizations have begun applying artificial intelligence (AI) and machine learning tools to automate various parts of the news production workflow (Carlson, 2015; Graefe, 2016). In practice, newsrooms deploy AI to handle routine and data-heavy tasks, enabling journalists to focus on analysis and investigative work (McFarland, 2015). This section surveys two broad areas: AI-driven workflow automation with software agents, and the integration of data pipelines and AI systems into news operations.

4.1.1 Workflow Automation with AI Agents

Modern newsrooms are increasingly turning to AI-powered agents and bots to handle repetitive tasks and broaden their coverage. These tools are capable of generating stories from structured data, managing content distribution schedules, and even tracking media feeds for breaking news. A notable example is the Associated Press (AP), which uses the Wordsmith platform from Automated Insights to automatically generate thousands of earnings reports from corporate financial data (McFarland, 2015).

In recent years, these systems have become even more advanced. In 2024, the AP rolled out its new Storytelling platform, which incorporates AI to help journalists plan, write, and publish content across different platforms all from a unified dashboard. This tool includes AI agents that suggest the best distribution strategies and help repurpose content assets, all while keeping human editors in the loop to ensure transparency and quality control. Likewise, The Independent uses a system called Bulletin, powered by Google's Gemini LLM, to create article summaries that editors then review and approve before they go live. These tools point to a growing trend: humans and AI are now working closely together in the editorial process.

Beyond content creation, AI bots are helping with editorial support behind the scenes. They can automatically tag articles with metadata-like topics, locations, or names making it easier to organize and retrieve content (Carlson, 2015). Some platforms are going further. Aftonbladet's Spånaren, for instance, uses past articles to recommend follow-up stories. AI agents also handle social media posts,

adjust content formats for mobile or multilingual readers, and power chatbots that respond to reader questions. The Washington Posts Climate Answers bot, launched in 2024, is one such example. It helps explain climate-related issues by drawing on the outlets past reporting.

Some current applications include:

- Generating templated reports from structured data (e.g., sports box scores, financial earnings) (McFarland, 2015).
- Tagging and organizing content (topics, keywords, sentiment) to streamline editorial workflows (Carlson, 2015).
- Scheduling and posting news updates on social media (maximizing audience engagement).
- Monitoring breaking news sources (newswires, social media) to alert editors of emerging events (Zhang et al., 2019).
- Summarizing complex topics with AI agents (e.g., Bulletin at The Independent).
- Recommending follow-up coverage based on existing article archives (e.g., Aftonbladets Spånaren).
- Facilitating multi-platform content planning with centralized AI-enhanced dashboards (e.g., AP Storytelling platform).

These agent-driven systems are already being used across the industry. Reuters News Tracer flags potential stories on X in real time, while Cleveland.com uses generative AI to draft basic community updates freeing up journalists for more investigative assignments. Meanwhile, Il Foglios experiment with a fully AI-generated edition showed both whats possible with automation and where it falls short. Ultimately, these technologies can dramatically boost newsroom productivity and reach, but human editorial oversight remains essential to uphold standards of accuracy and credibility.

4.1.2 Integration of Data Pipelines and AI Systems

Beyond standalone bots and agents, many newsrooms today are weaving AI and data analytics directly into their publishing pipelines from start to finish. These end-to-end systems manage everything from pulling in raw data to helping editors make decisions and pushing final stories out to audiences (K. Dörr, 2016; Graefe, 2016). Rather than being isolated tools, AI features are now embedded throughout the journalistic workflow. This kind of integration helps ensure that automation enhances every step of the process.

At the start of the pipeline, ingestion tools pull in a wide range of data. These sources can include structured feeds like financial databases, corporate filings, or sports scores, as well as unstructured content like tweets, RSS feeds, and open-data government portals. Platforms such as Apache Kafka and Apache Nifi help normalize and stream all this data in real time to the next stage. Many organizations also tap into web scraping tools (like Scrapy) and APIs (like X API or NewsAPI) to enrich their inputs even further.

Next comes processing, where machine learning models begin to interpret and enhance the content.

NLP tools often built on transformer models like BERT or RoBERTa are used to categorize articles by topic or tone, and to identify named entities like people or places (Devlin et al., 2018). Visual content goes through computer vision systems such as Google Vision or AWS Rekognition for object recognition, facial analysis, and scene interpretation. Some systems even support multilingual workflows, translating or summarizing content across languages. It's also here that editorial metadata like geographic focus, bias, or sentiment is often added to help guide future content curation or personalization.

A few real-world platforms show how this works in practice. The Washington Post's ARC system, for instance, uses cloud-based NLP tools to auto-tag and cluster stories, making editorial workflows much faster. In Europe, Schibsted has built its own AI infrastructure to process massive volumes of content for categorization, recommendation, and ad targeting. They use orchestration tools like Airflow and serve their models with TensorFlow Serving. Meanwhile, the BBC's Juicer system uses a microservices setup to handle ingestion and entity tagging, letting staff across different departments search and reuse enriched content easily.

The big advantage of this kind of integrated pipeline is that it creates a flexible, scalable editorial system where AI supports humans at every turn. Drafts can be auto-generated from structured data, articles can be tagged automatically, and potential stories can be flagged based on algorithmic scoring all before a human editor even steps in. While editorial judgment and ethical responsibility still rest with people, these systems help boost speed, accuracy, and output across the board.

4.2 Generative AI for News Content

Advances in generative models have created new possibilities for news writing and editing. This section focuses on AI that produces text: both rewriting existing text and generating novel content. Generative AI can create full articles, convert bullet-point data into narratives, or draft summaries and headlines.

4.2.1 Text Generation and Rewriting Models

Text generation generative AI varies from dated rule-based technology to state-of-the-art neural models. On one end, template systems are still commonly used in newsrooms, particularly for generating reports from formatted data, such as financial reports or sports summaries. These are pre-defined template-based systems in which data inputs (e.g., figures, dates, facts) are filled into pre-defined text structures. Though widely used in certain types of automated reporting (e.g., Automated Insights' Wordsmith, Narrative Science's Quill), template-based models are inherently rigid. They require extensive amounts of manual template design and maintenance, limiting accommodation for more flexible or creative content generation.

At the other end of the model spectrum, large pre-trained language models (LLMs) such as GPT-4, Meta's LLaMa, or Anthropic's Claude have revolutionized the ability to generate news content. These models are transformer-based and trained on enormous text corpora, allowing them to generate highly coherent and contextually dense text from very little input. In contrast to template systems, LLMs have the ability to learn to work with different writing styles and create new content without requiring

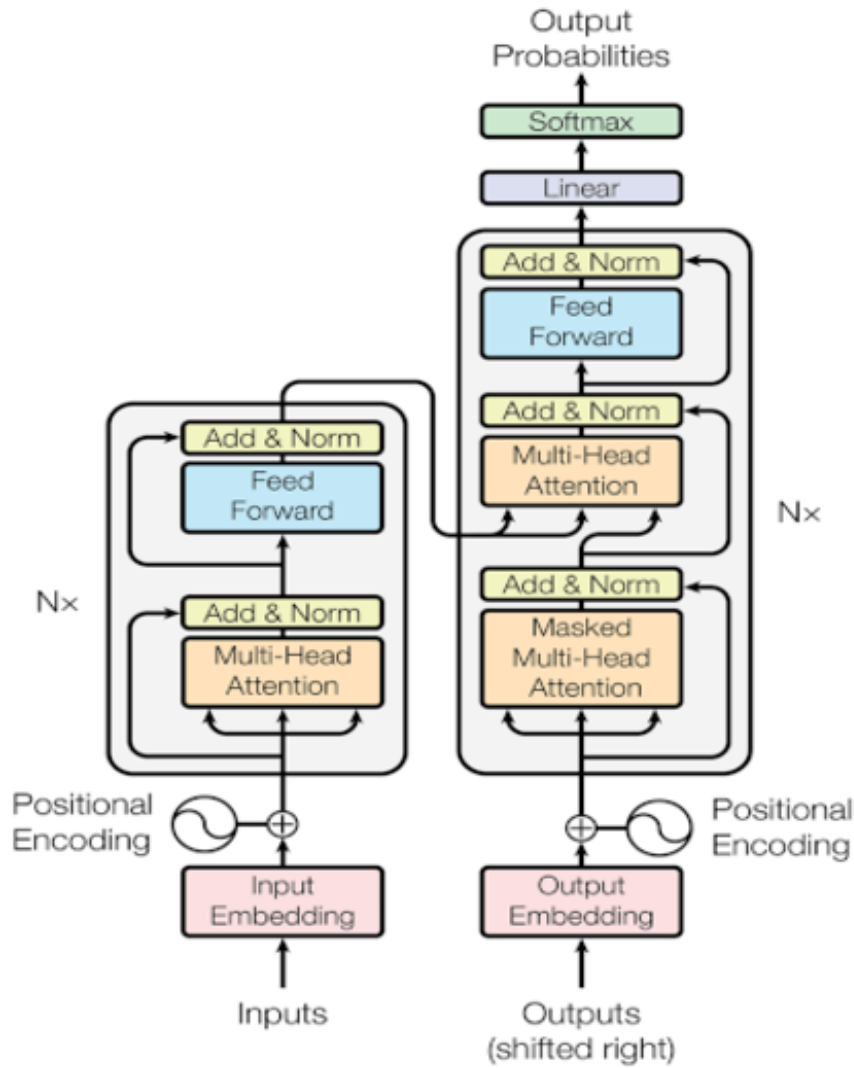


Figure 4.1: BERT Model Architecture

detailed manual specification. Fine-tuning LLMs on domain data (e.g., sports news, political rhetoric, or accounting reports) allows them to produce content in niche domains.

For instance, the ability of GPT-4 to generate rich stories from structured data is a breakthrough in dynamic content generation. This capability allows news agencies to mechanize writing functions like generating sports summaries or summarizing revenue reports in high fluency. Recent research, for example, Brown et al. (2020), has indicated that LLMs like GPT-3 can produce coherent sports recaps or short news pieces based on structured input (e.g., game statistics, market data) (Brown et al., 2020). They can also be further fine-tuned to adhere to specific editorial guidelines or styles, making them versatile content generation tools.

This notwithstanding, deployment of LLMs in newsrooms is often hybridized. Some use template-based systems for those tasks with significant amounts of required structure and standardization (such as election results or market summaries) but apply LLMs to fairly loose and creative tasks such as analysis pieces, interviews, or cover stories. Bloomberg, for example, has taken GPT-based tools to produce market summaries automatically, and the drafts are then used as starting points for editing and human fine-tuning (Lewis et al., 2020). This combination ensures that newsrooms can take advantage of the efficiency of AI without sacrificing editorial quality.

However, the use of LLMs for unsupervised writing in newsrooms is still relatively rare. In fact, AI-created drafts require a lot of editing by humans, especially for complex or sensitive topics. While LLMs like GPT-4 can produce syntactically correct text, human intervention still plays a role to ensure factual accuracy, delicacy, and strict compliance with ethical standards.

Recent advancements have made the deployment of generative models more accessible, with tools like OpenAI's GPT-3 and GPT-4 being integrated into user-friendly platforms like Jasper AI and Writesonic. These tools allow journalists to quickly generate first drafts, headlines, and summaries, streamlining workflow processes. Additionally, newer models, such as Meta's LLaMa, focus on open-source access, providing smaller news organizations with the ability to implement advanced generative AI without incurring the high costs associated with proprietary models.

In short, while LLMs are the cutting edge of generative AI in journalism, the direction in newsrooms these days is a mixed model where AI supports, not replaces, human journalists. Through automating the creation of routine content and suggesting editorial improvements, AI tools enable newsrooms to be more productive and focus on more investigative or advanced reporting. However, the effort of human editors remains vital in maintaining the integrity and credibility of AI-created content.

4.2.2 Automated Headline and Summary Generation

Automated headline and summary generation is one of the newer and growing areas of application in the use of AI technologies in newsrooms. As a component of the general trend to limit the effort that goes into content creation, AI systems are increasingly being used to generate brief, catchy short-form content, such as headlines, leads, and summaries. These systems use sophisticated machine learning models both to condense long articles into readable form, as well as to optimize them for viewing.

Beyond Rule-Based Approaches Headline generation used to be founded on naive rule-based systems or extraction of key phrases from articles. However, existing approaches leverage neural networks and massive pre-trained language models (LLMs) that are capable of capturing the context and nuances of a story. The nature of generated headlines has completely transformed with the shift from rule-based systems to transformer-based architectures such as BERT and GPT. They use contextual embeddings to identify the most salient facts of an article and generate headlines that are not only relevant but also engaging for readers.

Recent headline generation innovations focus on the ability to optimize for a specific outcome, e.g., click-through rate (CTR) or reader sentiment. AI systems now provide newsrooms with a list of potential headlines, which can be scored and ranked based on predicted engagement metrics. For instance, Hugging Face's Transformers library and GPT-4 are some of the platforms that have been fine-tuned for generating optimized headlines, with some systems specifically designed to predict audience engagement based on historical data (Xu and Lan, 2020).

This pattern is shown in The Washington Post's use of Heliograf (The Washington Post, 2020), an AI system that not only generates automated headlines but also rewrites them in real-time based on reader response. Heliograf, originally introduced for producing election results, now generates headlines

for a variety of topics, learning to craft permutations that maximize engagement across platforms (Vega, 2021). This indicates a growing shift to AI that learns and adapts based on reader interests.

As the technology for summarization improves, we are witnessing the advent of advanced abstractive summarization systems that generate human-like summaries, unlike the traditional extractive summarization that simply extracts key sentences. SOTA transformer-based models like T5, BART, and PEGASUS are at the forefront of abstractive summarization for newsrooms. They are trained on massive corpora of news articles and can generate summaries that are not only fluent but also dense contextually, conveying the key points of an article in a summary efficiently.

For example, PEGASUS developed by Google Research was specifically geared towards abstractive summarization, building on a pretraining strategy involving masking portions of input articles and having the model generate the masked sections (Lewis et al., 2020). As a result, PEGASUS can produce high-quality summaries of lengthy news articles that are accurate and concise, even when summarizing content that is made up of multiple paragraphs or contains nuanced subject matter.

News organizations are increasingly adopting these advanced summarization models to create a wide range of short-form content, from news in 60 seconds videos to social media summaries. One example of this is the use of OpenAI's GPT-4 to automatically generate summaries for breaking news, where the AI condenses long news articles into succinct summaries for platforms like X and LinkedIn. In some cases, such AI-generated summaries are directly fed into content management systems (CMS) to speed up the publishing process, allowing newsrooms to catch up with the competitive online news environment.

Most notable among developments in headline generation has been the inclusion of engagement-driven algorithms that optimize headlines for higher user engagement. These systems use machine learning to analyze previous data in order to predict which headline formats will engage a particular audience most. Emotional tone, relevance, and readability are all evaluated to generate headlines that are most likely to generate clicks or social media shares.

The New York Times and other large publications have begun leveraging AI-powered systems to A/B test versions of headlines to assess their potential performance. Such systems continually refine their predictions based on real-time data feedback loops, which allows for continuous optimization of headlines in relation to engagement trends (Su and Xu, 2019). This capacity is also reinforced by Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) models that can better understand the sequential nature of headlines and audience response over time.

While there have been incredible strides in automatic headline and summary generation, such systems still falter, particularly in terms of fact checking and journalistic ethics. AI systems, particularly with abstractive summarization, sometimes generate summaries that omit crucial information or even include factual inaccuracies. A headline generated by an AI system, for example, might highlight an overly sensational or dramatic aspect of the story, which would be inaccurate if not put through editorial oversight.

Moreover, AI-generated headlines and summaries must adhere to the ethical guidelines of journalism

so that the tone is appropriate and the content is not a distortion of facts. While it is possible to achieve high-quality outputs using tools like GPT-4 and PEGASUS, such models must be subjected to rigorous fact-checking and human review to avoid potential issues like biased language or flawed framing.

While AI-generated content continues to gain momentum, the role of human editors remains crucial. While AI software can speed up the editing process and help journalists ramp up content production, editorial staff continue to have the final decision on what goes to print, ensuring that AI products pass the publication's tests of accuracy and ethics.

In conclusion, automatic headline and summary generation is a rapidly evolving field, with recent advances driven by deep learning and transformer-based models. AI deployment in newsrooms has allowed content creation to become more efficient, scaling content production for organizations at the short-form end while still allowing for some level of personalization and optimization for reader engagement. Yet despite enormous technological breakthroughs, the need for human supervision remains necessary to ensure the integrity, truthfulness, and ethical considerations of the AI-generated content.

4.3 Machine Learning for Editorial Assistance

Beyond content generation, machine learning plays a growing role in supporting editorial workflows by providing analytical tools that assist rather than automate decision-making. These systems help editors classify stories, anticipate audience interest, and detect emerging trends. The emphasis is not on writing text but on augmenting editorial oversight through data-driven insights. This section reviews two key applications: content classification and predictive analytics.

4.3.1 Classification of News Content by Topic and Tone

Machine learning classifiers are often used in newsrooms to label articles by topic, sentiment, and urgency to enable better content organization and delivery. Topic classifiers, which are often built using supervised learning algorithms such as Support Vector Machines (SVMs) or fine-tuned transformers (e.g., BERT), can label an article with multiple labels through headline and body text analysis (Choi, 2018). Labels can capture thematic categories (e.g., climate, politics, health) or more detailed beats based on editorial needs.

Tone and sentiment detectors analyze subjective news elements, whether content is neutral, positive, critical, or charged. This is particularly relevant in editorial balance and in routing sensitive news. For example, models that detect urgency identify rapidly changing news requiring moment-to-moment editorial correction. Some CMS-based tools also detect toxic or libelous language, which prompts editors to review highlighted passages before publication.

These classification frameworks are often incorporated into newsroom content management systems (CMS). When a reporter submits a draft of an article, the system automatically labels it with metadata such as categories, keywords, sentiment scores, or named entities. This metadata facilitates sophisticated querying—editors can, for example, look for recent negative-tuned coverage of economic policy or find

coverage that mentions a specific public figure by name. This automation relieves the tedium of hand tagging and facilitates discoverability and content tracing over broad archives.

4.3.2 Predictive Analytics for Newsworthiness and Audience Engagement

Predictive analytics applies machine learning to forecast how stories might engage with consumers based on historical data and content features. Predictive models forecast possible engagement metrics such as page views, dwell time, shares, or conversion probability of a subscription. Some early work demonstrated that article popularity on social media could be predicted using features like topic, named entities, or publication time (Bandari et al., 2012). Subsequent innovations have introduced deep learning architectures and real-time feedback loops to continue refining such predictions (Su and Xu, 2019).

News organizations now use predictive models for both planning and real-time editorial decision-making. Internally trained models can rate draft articles and rank them based on estimated reach or reader engagement, assisting editors in determining what to push or revise. These models are able to consider a mix of traits: textual (for example, sophistication, tone), contextual (e.g., timing, competition), and behavioral (e.g., reader history). Multiple headline variations can also be offered with engagement ratings by some systems so editors can select the best one to pursue.

Machine learning is also used in following external signals. Tools such as *Dataminr*, *Google Trends*, or *NewsWhip* utilize real-time data mining and anomaly detection to surface trending stories in a rush and alert editors of viral or breaking news. For instance, *Reuters' News Tracer* system analyzes the velocity and credibility of tweets to detect early signals of emerging stories, even ahead of traditional wire services Zhang et al., 2019. These trend-spotting systems are underpinned by semi-supervised and unsupervised models that group and rank information sources on novelty, volume, and credibility.

There is a new focus on editorial analytics being *explainability*: ensuring predictive systems deliver transparent explanations of recommendations. Feature attribution and attention mapping are some of the methods employed to allow editors to understand why some stories are ranked higher, gaining trust in AI outputs and enabling more informed editorial decisions (Kumar et al., 2022).

Though progress has been made, predictive systems remain editorial guidance tools. They are meant to assist, not replace, editorial judgment—making probabilistic recommendations, not certain conclusions. Editors still apply their judgment in judging newsworthiness, maintaining ethics, and choosing publication. Machine learning facilitates this process by exposing patterns and clues that would be difficult to identify by hand, especially in the high-velocity digital news world.

4.4 Human-in-the-Loop vs. Full Automation

One of the key questions in AI-powered journalism is how much control humans should have over the content pipeline. Most modern newsrooms lean toward a *human-in-the-loop* (HITL) approach, where AI tools assist journalists but don't take over. In this setup, systems might generate a draft, headline, or summary, but a human editor always steps in to review, tweak, and sign off before anything goes live (Weedon, 2021). This layer of editorial oversight helps maintain accuracy, ethical standards, and

the outlets unique voice. As Marconi (2019) points out, top media companies usually require human review especially when the content touches on sensitive legal or reputational ground (Marconi, 2019).

Weve already seen this model in action. The Associated Press uses natural language generation (NLG) to crank out earnings reports from structured financial data but editors still review each piece before its published. Bloomberg follows a similar method, automating market summaries but keeping humans in the loop to make sure everything is factually and contextually sound. Outlets like The Guardian and the BBC have also tested AI-assisted writing tools, but journalists remain firmly in charge when it comes to tone, accuracy, and final approval (Weedon, 2021).

On the flip side, theres the *full automation* route publishing without any human review. This is usually reserved for content thats low-risk and high-volume, like sports scores, weather updates, or stock tickers (Sawers, 2014). In these cases, the data is clean, the structure is predictable, and errors are easy to catch. For example, RADAR in the UK has used structured feeds to publish thousands of local news briefs every month, all without human editing.

That said, fully automated publishing isnt the go-to for complex or sensitive stories. When language gets tricky or topics are controversial, human oversight is still essential. There are ongoing concerns about AI hallucinating facts, embedding bias, or mishandling sources (Graefe, 2016). Even with powerful language models, the need for responsible editorial judgment hasnt gone away.

A hybrid model is proving to be the sweet spot for many outlets. It gives them the scale and speed of automation while preserving trust and journalistic standards. These flexible workflows let teams switch between full automation and HITL depending on the topic, urgency, or risk level making room for smarter, more effective human-AI collaboration.

Chapter 5

Designing the Journalist Journey: Tracing a News Article in an AI-Powered Newsroom

5.1 Introduction

As we have seen throughout the overall project, the journalist's work has been deeply modified by the advent of artificial intelligence and technologies in its widest definition. In this chapter we will illustrate the roadmap a news follows, from its very first embryonic state to its publication. We will mention the main technologies used, their functioning and their purpose.

The exemplification of these processes will be around a not announced protest event occurred in the city of Dondure, in the fictional country of Estania, in April 2025, where a spontaneous demonstration erupted after a controversial AI law was approved overnight by the national parliament. The protest, initially unnoticed by mainstream media, was picked up through social media buzz and foreign journalist networks before becoming a breaking news item for international press agencies.

Our Journalist is in charge of political news and this event perfectly aligns with its role.

5.2 News Gathering

The day of the journalist of course starts with a good coffee. Its day cannot start without it. After this, it usually check for emails and possible updated on its previous works. Nothing too difficult, AI here is still not implemented.

In the course of becoming constantly up to date about global developments, journalists rely on a multitude of information resources implemented and maintained by their newsroom newsroom that we can refer to as *The Flux*. The Flux is an ancient institution news agency operating under the optimum of global standards. Thanks to its strong global presence and reputation, it has made official arrangements with quite a few other top news agencies around the world from South America to China, from Spain to the United States. These arrangements deliver real-time flows of information in multiple forms

like API integration, RSS feeds, wire services with structured outputs, and secure editorial feeds. This implies that the agency can offer comprehensive and responsive coverage of breaking news on every continent and filter and authenticate inputs through both automatic and editorial controls.

Among the incoming sources handled by The Flux, one of the most structured and constant streams comes from formal agreements with international news partners. These agencies provide their content in machine-readable formats, most commonly XML (Listing 5.1), with each item clearly labelled with metadata such as source identifier, language, publication time, and topic tags.

Listing 5.1: Example of raw XML feed from Estania

```
<newsItem source="EstaniaPress" language="eng" \n
  timestamp="2025-04-12T08:13:00Z">
  <title>Unusual banner spotted near Estanian Parliament
</title>
  <body>
    A cloth banner was seen early this morning
    hanging from a pedestrian bridge near the
    Parliament building in Dondure. The sign,
    written in spray paint, appears to reference recent
    discussions on AI regulation.
    Authorities have not commented.
  </body>
  <tags>
    <tag>politics</tag>
    <tag>AI</tag>
  </tags>
</newsItem>
```

Upon opening the platform where news arrives, the journalist receives only those stories tagged with topics relevant to their area of focus. For our journalist, who specializes in politics, they will primarily receive news tagged with terms such as "politics", "referendum", "elections", or "parliament". At first glance, nothing seems particularly noteworthy a few political statements, parliamentarians debating over funds, and a banner in Estania. These are everyday occurrences that happen frequently and in significant numbers. Journalists often struggle to identify which stories are truly relevant amid the flood of information they are presented with (3.2.1). As a result, it is often difficult to pinpoint the right topic to highlight at first glance.

As we have mentioned in Section 3.3.2, in the today's news environment, newsroom do not only have to compete with traditional competitors, but also with less formal and more instant actors, which mainly work in the social media world. Specifically, posts, tweets or hashtag trends are fundamental sources to get informed about events.

To face with such phenomenon, The Flux has implemented different strategies. Thanks to the usage of APIs some of the main social media are covered and constantly monitored in order to identify, which

are the main trends. Listed the techniques used for different social medias:

- **Reddit:** Reddit is a wide open social where public posts and subreddits are used by followers to discuss about specific topics. It is a powerful source for identifying trends. Specifically, the **PRAW Library** in Python is a powerful tool for interacting with Reddit. For each specific area of the newsroom, there are defined settings for the usage of the library and display of information for the user. Subreddits such as "r/politics" or "r/PoliticalDiscussion" are everyday monitored. It is possible to filter for most trendy discussion or comments based on their score (upvotes).
- **Instagram:** Instagram is a visually-driven social media platform that, while more restricted than others, still offers valuable signals for trend detection especially via hashtags. In a newsroom context, public data from business or creator accounts can be accessed through the official **Instagram Graph API**, part of the broader Meta API ecosystem. This access requires prior authentication and connection to a verified Facebook Page. For integration into automated monitoring workflows, Python scripts using libraries such as **requests** are configured to interact with endpoints like `ig_hashtag_search` and `recent_media`, allowing retrieval of metadata such as captions, timestamps, and media links. Dedicated internal modules allow journalists to specify hashtags of interest (e.g., `#politics`), and the system periodically fetches new content for further NLP analysis or human review.
- **X:** X is maybe the most useful source for newsrooms in order to identify trends, thanks to its wide usage and hashtags mentioning. X has an official API (API v2) which allows to interact with X, collecting metrics such as retweets, likes, or replies count. Our newsrooms uses the python library **requests**, making HTTPS requests, for specific hashtags or accounts posting.

The Flux uses various techniques to identify key topics and trends within the massive amount of data collected. By applying Natural Language Processing (NLP) techniques such as topic detection, sentiment analysis, and anomaly detection, the system can automatically flag relevant posts or threads.

Once the data is processed, it is presented to journalists through a visualization dashboard. The dashboard aims to simplify the massive datasets, making them actionable for journalists. It includes several key features:

- **Trending Topics:** A dynamic list of hashtags or keywords that have experienced significant growth over a short period.
- **Geolocation Data:** A map visualizing where specific hashtags or mentions are being discussed most, providing insight into the geographical spread of a topic.
- **Sentiment Analysis:** A summary of the overall sentiment surrounding a trend or hashtag, indicating whether the conversation is predominantly positive, negative, or neutral.
- **Time-based Trends:** Graphs that track the rise and fall in mentions of certain hashtags over time, providing context on how the trend evolves.

Listing 5.2: Trending Posts from r/politics

```
import praw

# Get top hot posts from r/politics
```

```

subreddit = reddit.subreddit('politics')
for post in subreddit.hot(limit=10): # Hot posts
    print(post.title, post.score, post.url)

```

Breaking News

When accessing the visualization dashboard the journalist sees a particular trend in hashtags related to politics discussions on reddit. A new hashtag **#EstaniaProtest**, is flagged by the system and appears to be trendy on social networks. The journalist has to inspect the realness of such, and decides to check the plots automatically made by the interactive dashboard.

What they immediately face is a clear trend for an apparent Estania Protest. Checking a simple time series plot (Figure 5.1), there is evidence of an occurring event.

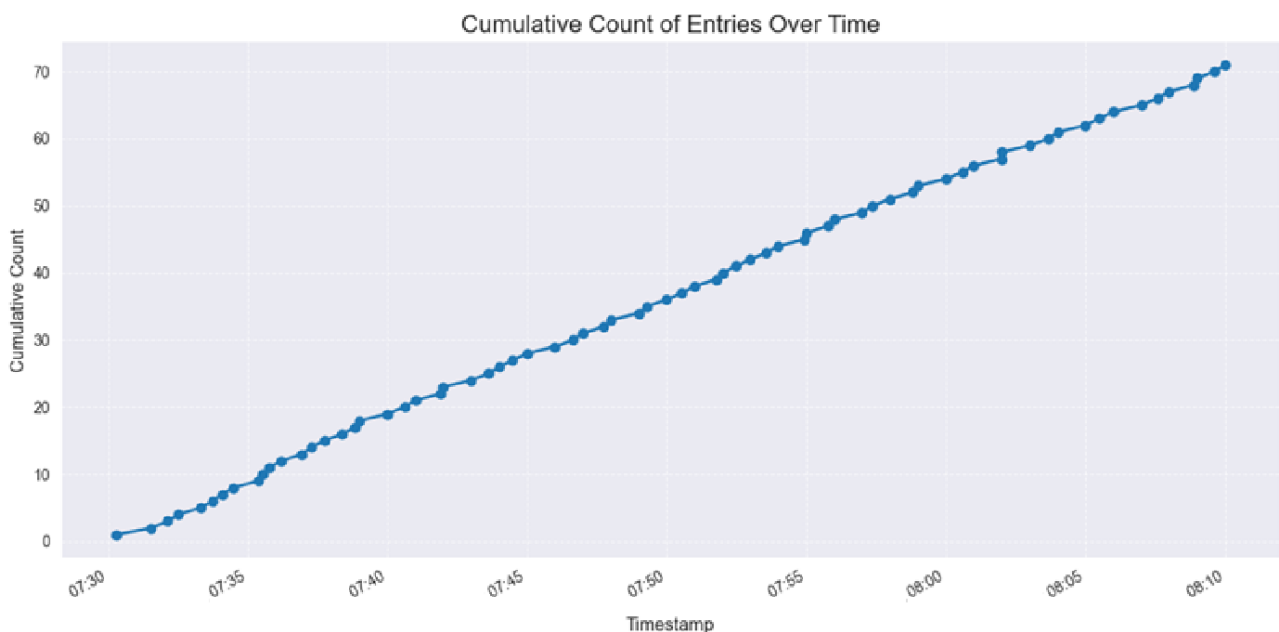


Figure 5.1: Time Series Plot for **#EstaniaProtest**

The journalist then checks if these posts exist and which is their content, by simply checking the posts that present the hashtag.

Listing 5.3: Simulated X XML response for **#EstaniaProtest**

```

<response>
  <data>
    <post>
      <id>1</id>
      <text>Protesters have gathered in Dondure,
      demanding justice
      for the Estania case. #EstaniaProtest</text>
      <created_at>2025-04-12T08:22:00Z</created_at>
    </post>
    <post>

```

```

    <id>2</id>
    <text>Reports of increasing tensions in
    downtown Dondure.
    People are blocking streets. #EstaniaProtest</text>
    <created_at>2025-04-12T08:45:00Z</created_at>
  </post>
</data>
</response>

```

Listing 5.4: Simulated Instagram XML response for #EstaniaProtest

```

<response>
  <data>
    <post>
      <id>180027363415</id>
      <caption>A strange banner seen in downtown Dondure today.
      #EstaniaProtest
      </caption>
      <media_url>https://instagram.com/media/xyz.jpg</media_url>
      <timestamp>2025-04-12T08:21:00+0000</timestamp>
    </post>
    <post>
      <id>180028874289</id>
      <caption>People start gathering in Dondure's main square.
      Tensions rising? #EstaniaProtest #politics </caption>
      <media_url>https://instagram.com/media/abc.jpg</media_url>
      <timestamp>2025-04-12T08:46:00+0000</timestamp>
    </post>
  </data>
</response>

```

The emerging picture suggests a protest forming in Dondure, possibly linked to recent developments in national AI legislation.

In breaking news scenarios, the journalist must react quickly. After a swift verification of post authenticity, a preliminary article is drafted by the journalist for The Flux website:

BREAKING Early Signs of Unrest in Estania's Capital Over AI Regulation Bill

Dondure, Estania April 12, 2025 08:47 AM (local time)

Early this morning, signs of unrest have begun to emerge in downtown Dondure, where scattered reports, photos, and videos on social media suggest the formation of a spontaneous protest against

Estanias newly proposed AI regulation law.

Initial images circulating on Instagram and X show protest banners and small groups forming near the Parliament building. Posts with hashtags like `#EstaniaProtest`, `#Allaw`, and `#Politics` began trending around 08:00 AM, suggesting a rapid mobilization.

So far, no official statement has been released by local authorities. However, several witnesses have posted footage indicating that a group of citizens is calling for the government to withdraw the proposed law, which critics say could limit digital rights and increase surveillance.

The Flux newsroom is actively monitoring the situation as it develops and will provide further updates as more verified information becomes available.

5.3 News Processing

After the first initial news publication, the journalist has to create a more informative and enriched report, which will be published on the next day on "The Flux" newspaper. The work of getting additional information would take a lot of time with traditional methods, but thanks to the informational retrieval technologies implemented by the newsroom, this work sounds way faster.

The specific technologies for this phase are:

- An SQL Database for Context Information Retrieval
- A Chatbot for interacting with database
- Natural Language Processing tools

First of all, the journalist looks for contextual information in order to better define the reasons and the possible actors involved in the event.

This work is supported by a internal database of the newsroom, containing all the articles published by "The Flux". Of course this process cannot be done manually by looking to all the articles of a said section.

The Flux has implemented a very good working AI assistant chatbot.

The AI assistant is based on a Retrieval-Augmented Generation (RAG) architecture, which allows the model to retrieve relevant information from a large corpus before generating a response. The process involves two main phases: **retrieval** and **generation**.

First, all documents in the newsroom database (past articles, reports, statements, etc.) are pre-processed and transformed into dense vector representations (also called *embeddings*). Specifically our newsroom uses the `multilingual-e5-large-instruct` due to its multilingual capabilities, strong performance, and instruction-

tuned nature, which is beneficial for information retrieval tasks Enevoldsen et al., 2025. These embeddings capture the semantic meaning of the text and are stored in a vector index, typically implemented with a similarity search library like FAISS.

When the journalist asks a question or makes a query (e.g., “*Have there been other protests in Estonia related to AI laws?*”), the assistant encodes the question into an embedding in the same vector space. Then, a similarity search is performed to retrieve the top-k most semantically similar documents from the database.

To enhance precision, documents can be clustered during preprocessing using unsupervised techniques such as K-Means or HDBSCAN, allowing the system to focus retrieval on thematically coherent subsets (e.g., protests, AI laws, political opposition).

Finally, the retrieved texts are passed to a generative language model that conditions its answer on this external information. Our newsroom LLM is GPT-4o mini, which can guarantee good performances with a competitive price (OpenAI pricing website), balancing the trade-off cost-quality. This architecture allows the assistant to remain grounded in verifiable data, providing accurate, context-aware support to the journalist.

	Rank (↓)	Average Across		Average per Category								
Model (↓)	Borda Count	All	Category	Btxt	Pr	Clf	Clf	STS	Rtrvl	M. Clf	Clust	Rnkr
MTEB(Multilingual)												
Number of datasets (→)	(132)	(132)	(132)	(13)	(11)	(43)	(16)	(18)	(5)	(17)	(6)	
multilingual-e5-large-instruct	1 (1375)	63.2	62.1	80.1	80.9	64.9	76.8	57.1	22.9	51.5	62.6	
GritLM-7B	2 (1258)	60.9	60.1	70.5	79.9	61.8	73.3	58.3	22.8	50.5	63.8	
e5-mistral-7b-instruct	3 (1233)	60.3	59.9	70.6	81.1	60.3	74.0	55.8	22.2	51.4	63.8	
multilingual-e5-large	4 (1109)	58.6	58.2	71.7	79.0	59.9	73.5	54.1	21.3	42.9	62.8	
multilingual-e5-base	5 (944)	57.0	56.5	69.4	77.2	58.2	71.4	52.7	20.2	42.7	60.2	
multilingual-mpnet-base	6 (830)	52.0	51.1	52.1	81.2	55.1	69.7	39.8	16.4	41.1	53.4	
multilingual-e5-small	7 (784)	55.5	55.2	67.5	76.3	56.5	70.4	49.3	19.1	41.7	60.4	
LaBSE	8 (719)	52.1	51.9	76.4	76.0	54.6	65.3	33.2	20.1	39.2	50.2	
multilingual-MiniLM-L12	9 (603)	48.8	48.0	44.6	79.0	51.7	66.6	36.6	14.9	39.3	51.0	
all-mpnet-base	10 (526)	42.5	41.1	21.2	70.9	47.0	57.6	32.8	16.3	40.8	42.2	
all-MiniLM-L12	11 (490)	42.2	40.9	22.9	71.7	46.8	57.2	32.5	14.6	36.8	44.3	
all-MiniLM-L6	12 (418)	41.4	39.9	20.1	71.2	46.2	56.1	32.5	15.1	38.0	40.3	

Table 5.1: MTEB (Multilingual) Benchmark Results
Source: Enevoldsen et al., 2025

All the articles provided by the assistant response bring along some metadata, containing additional information about the article. These metadata are a pure output of Natural Language Techniques applied. These metadata bring information like:

- Entities Extracted from the Text. Particularly:
 - Text of the Entity
 - Category of the Entity (a person, an organization)
 - A description of the entity

- Sentiment Analysis score
- Topics treated in the article

These metadata do not only allow for a better understanding of the article, but also for the clustering and speed of researches the journalist can perform.

5.3.1 Named Entity Recognition

The NER is performed again using the GPT-4o mini model, for the previously mentioned reasons. Through a prompt the model is charged of extracting entities from the text which pertain to one of the following categories:

- PER: Persons (John Smith, Kodak Black)
- ORG: Organizations (FIAT, Rolex)
- LOC: Locations (France, Palestine)
- MISC: Miscellaneous (Spanish, Marxism)

Once an entity is recognized by the model, a brief description of the entity based on the context of the article provided is inferred and returned in a JSON format.

```
{
  "text": "Dondure",
  "category": LOC,
  "description": "Capital of Estania"
}
```

5.3.2 Sentiment Analysis

Sentiment Analysis is used to understand the emotional tone of the articles retrieved, allowing the journalist to better assess public opinion or media tone towards specific topics or events. Unlike the NER task, sentiment analysis is handled by a fine-tuned open-source transformer model: **RoBERTa-base**, specifically trained for sentiment classification.

RoBERTa was chosen for its excellent balance between accuracy, performance, and cost-efficiency. As an open-source model developed by Facebook AI, it is freely available and can be integrated into internal systems without dependency on proprietary APIs. Its transformer-based architecture allows for contextual understanding of language and has been proven to outperform other models on standard sentiment datasets (e.g., SST-2, IMDb), reaching accuracy levels above 93%.

Once a given article is retrieved, the model processes the content and classifies the sentiment as either **positive**, **negative**, or **neutral**. It outputs both the classification and a confidence score. Additionally, an internal post-processing module generates a short justification for the sentiment label using keyword patterning or large language models (if needed).

An example of the output structure could be:

- `sentiment_label`: *negative*
- `confidence_score`: *0.82*

This information allows the newsroom to perform large-scale aggregations of public sentiment over time or across topics (e.g., AI laws, political decisions). Furthermore, sentiment labels contribute to filtering or ranking articles during the exploration process, helping the journalist quickly identify polarized or emotionally charged content.

Model	Accuracy(%)	Precision(%)	Recall(%)	F1-Score(%)
BERT (base-uncased)	91.8	91.8	91.8	91.8
RoBERTa (base-uncased)	93.4	93.5	93.5	93.3
XLNet (base-uncased)	92.5	92.5	92.5	92.5
ALBERT (base-v2-uncased)	91.4	91.4	91.4	91.4

Table 5.2: Performance comparison of transformer-based models on the SST-2 dataset.
Source: Jiang, 2020

5.3.3 Topic Modelling

Topic modelling is another key NLP technique adopted by The Flux newsroom to support journalists in understanding the underlying themes present in historical or trending articles. Given the complexity and volume of content available, this tool helps cluster documents by latent semantic structures, enabling efficient navigation through related material.

For this task, The Flux leverages the Latent Dirichlet Allocation (LDA) algorithm – a probabilistic generative model that assumes each document is a mixture of topics, and each topic is a distribution over words. LDA is particularly useful for newsroom applications due to its interpretability and effectiveness in extracting human-readable themes from large corpora.

Each article is preprocessed through standard NLP techniques (tokenization, stopword removal, lemmatization), and then passed to the LDA model. The model outputs:

- A distribution of topics per article (e.g., 60% Topic A, 40% Topic B)
- The top keywords associated with each topic
- (Optional) A label or inferred category for each topic

These results are stored as metadata for each article and used in the newsroom platform to allow:

- Fast clustering of related articles
- Recommendation of background readings when breaking news emerges
- Filtering of the database by conceptual theme rather than keyword

Here an example of a possible output:

```
{
  "topic_id": 3,
```

```

    "keywords": ["protest", "ai", "law", "citizens", "rights", "regulation"],
    "weight": 0.65,
    "label": "Civil Unrest / AI Policy"
  },
  {
    "topic_id": 7,
    "keywords": ["parliament", "debate", "bill", "government", "vote"],
    "weight": 0.25,
    "label": "Legislation Process"
  }
}

```

The Journalist now has all these information in a easy to understand format and can now gain additional information about what the context is about and which are the entities involved.

The next step is then the generation of the content.

5.4 News Generation

The journalist has now gathered additional information and knowledge, which allows them to create a satisfactory and exhaustive report about the causes and developments of what happened.

Collecting both the information from the back context thanks to the chatbot, and merging the continuous social media tracking and analysis deriving from the web, the journalist has a clearer landscape of what happened and what is happening.

Through the RAG-based assistant, the journalist retrieves a 2023 article previously published by The Flux, reporting on a similar protest in Estania following a controversial surveillance law proposal. That protest had been partially organized by a decentralized group known as "FreeCode Estania", a civil liberties activist network advocating for algorithmic transparency. The archived article (ID: FLX-2023-145) described how the group used encrypted platforms and Reddit subforums to mobilize supporters.

Further exploration via entity tracking and co-occurrence analysis reveals that "FreeCode Estania" has been explicitly mentioned again in recent social media posts, some of which include direct claims of responsibility for the current protest (e.g., We did it in 2023, and well do it again until AI is fair. #EstaniaProtest).

Moreover, one politician in particular, Marlon Veznick, leader of the opposition Social Justice Party, is frequently referenced both in historical coverage and in present conversations. Several posts associate his name with supportive statements toward the protest. LDA-based topic modeling confirms that his partys public stance on AI regulation is gaining traction among protest narratives.

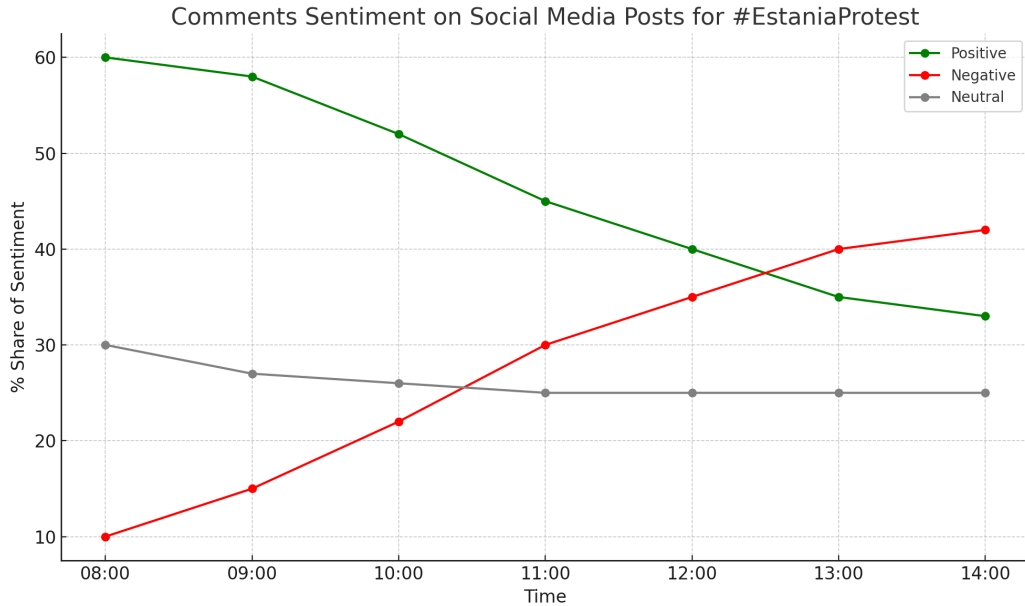


Figure 5.2: Sentiment Analysis from social media posts

From the continuous social media monitoring, sentiment analysis on posts tagged with **#EstaniaProtest** and **#Allaw** reveals a noticeable increase in negative sentiment, especially in replies and quoted posts. Hate-related vocabulary has spiked in messages during the day, according to the sentiment time series plot (Figure 5.2). Analysts warn that such a trend could signal the risk of further radicalization or potential escalation of unrest.

The journalist now has a much clearer and multi-perspective view of the storylinking past events, tracking actors involved, and understanding public reaction in real timeall of which enrich the journalists capacity to report not only on what is happening but why.

All this information is then passed on to the Generative AI model, which is able to produce a journalist-tone-like article, tailored to match the tone of The Flux. This is made possible by a fine-tuned large language model based on GPT-4, trained via Reinforcement Learning from Human Feedback (RLHF). In this process, human annotators rank multiple outputs generated from the same prompt, and a reward model is trained on these rankings. The RLHF process makes the model generate output in accordance with journalistic conventions through the optimization of human-typed preference outputs in terms of tone, shape, and factual consistency.

- Contextual embeddings: Combined metadata from social media (e.g., sentiment trends, extracted entities, timestamps) and the newsroom’s internal knowledge base are encoded into dense vectors using Sentence-BERT.
- Template-aware generation: Style-specific prompts are used to condition the generation toward specific formats (e.g., breaking news, analysis, summary).
- Post-editing filters: Generated text is then passed through a sanity and fact-checking layer that verifies entity consistency and flags hallucinations using a dedicated named entity validator and keyword consistency checker.

This generation framework allows The Flux journalists to receive draft articles tailored to their editorial style, yet fully editable for final publication.

5.4.1 Generated Article

Here an example a possible final article produced by the AI assistant.

ESTANIA Growing Tensions as Historic Protest Escalates Over AI Law Reform

Dondure, Estania April 12, 2025, 2:00 PM

Tensions continue to escalate in Dondure as the protest against the proposed AI regulation bill grows in size and intensity. Following early signs of unrest this morning, large crowds have now filled the city's central avenues, chanting against what they claim is an authoritarian overreach masked as AI ethics.

Social media monitoring platforms, leveraged by The Flux newsroom, show a clear rise in posts tagged with #EstaniaProtest, #AIlaw, and #NoToSurveillance, many of which now carry a marked increase in negative sentiment. An analysis of Instagram comment threads reveals a sharp drop in positive engagement, particularly following the government's announcement at 10:00 AM, which was met with widespread criticism.

Further historical context retrieved through The Flux AI assistant revealed that a similar wave of protests occurred in late 2023, organized by the now-active civic group Digital Liberties Front (DLF), known for its opposition to digital surveillance. Their official channel posted a statement this afternoon claiming responsibility for coordinating today's march and warned of further mobilizations if the bill is not withdrawn.

The role of opposition leader Marta Kelvik (Social Democratic Front) is also under scrutiny, as she has repeatedly shown support for the protestors' cause across multiple platforms.

With the situation evolving rapidly, and narratives online becoming increasingly polarized, concerns rise over the potential for clashes with security forces.

The Flux will continue to monitor developments and provide timely updates as more facts emerge.

5.4.2 Human-in-the-Loop Validation

Despite the capabilities of advanced Generative AI systems, *The Flux* maintains a fundamental principle: no AI-generated article is published without final human review. After the language model produces a draft based on retrieved information and enriched metadata (e.g., entities, sentiment, topic), our journalist intervenes to assess:

- **Factual Accuracy:** Verifying every claim against trusted sources or the newsroom’s database.
- **Narrative Coherence:** Ensuring the structure and flow of the article align with professional standards and the editorial tone.
- **Ethical Responsibility:** Evaluating bias, representation of individuals or groups, and potential societal impact of publishing.

This human-in-the-loop approach guarantees that while AI accelerates information processing and drafting, editorial responsibility remains firmly in human hands. This practice ensures transparency, accountability, and upholds journalistic integrity in the age of automation.

5.5 News Distribution

With the journalist having completed their final review and approved the AI-generated content, the process enters its final stages: those of automated publishing and personalized distribution.

5.5.1 Publishing

Before content is released publicly, an automated system performs a final layer of validation to ensure lexical coherence, stylistic uniformity, and tonal alignment with the outlet’s editorial standards. This is achieved through Natural Language Processing tools that analyze the syntactic and semantic patterns of the text and compare them with a corpus of past *The Flux* articles. Any inconsistencies or outliers are flagged for optional secondary review.

Once validated, the article is passed to an integrated Content Management System (CMS). This system, connected via secure APIs, handles:

- **Automated Scheduling:** Articles can be scheduled based on optimal traffic hours or editorial calendars.
- **Structured Metadata Injection:** The final article includes structured metadata extracted or inferred during the enrichment phase (e.g., NER, sentiment, topic modeling). This boosts SEO performance and enables semantic interoperability across platforms such as news aggregators, social media, or archive services.
- **Multichannel Adaptation:** The article is automatically formatted for different digital outlets: e.g., responsive layout for web, condensed version for mobile, or teaser snippets for social feeds.

5.5.2 Translating & Targeting

In order to maximize reach and relevance, *The Flux* applies a multilingual and audience-aware dissemination strategy.

Machine Translation Generative AI, based on the GPT-4o mini model is employed to translate articles into Italian, Spanish, and French. While the initial translations are automated, a layer of assisted post-editing is applied where linguistically-aware tools suggest improvements, and editors intervene only when necessary. This ensures rapid yet high-quality multilingual output.

Audience Targeting Distribution is fine-tuned according to:

- **User History:** Behavioural analytics identify users who have previously read about related topics (e.g., AI regulation, protests).
- **Geolocation:** Based on user location (e.g., Estania, Brussels), priority content is surfaced that aligns with local interest or impact.
- **Predictive Engagement Models:** Machine learning models, trained on historical click-through and reading time data, predict which headlines, formats, or themes are more likely to engage specific users.

Personalized Output The headline, preview snippet, or even the article’s introduction can be subtly adjusted based on the target distribution channel. For example, a more emotionally resonant headline may be selected for Instagram or newsletter placements, while an informational tone may be preserved for the homepage.

In short, the pipeline’s process from real-time event detection to cross-lingual publication is that of an extremely integrated AI-enabled ecosystem. The journalist remains crucial in steering quality and ethical judgment, but the heavy lifting of scaling, adapting, and targeting content is performed by an effective stack of language and recommendation models. This closes the loop of an enabled contemporary newsroom with language technologies.

Chapter 6

Discussions and Conclusions

6.1 Revisiting Research Questions

The research questions introduced in Chapter 1 provide a roadmap for this study. We summarize how each question has been addressed through literature review, case studies, system simulation, and critical analysis:

1. **RQ1: How is artificial intelligence currently implemented in news production workflows?** We found that AI is used in multiple stages of the news pipeline, including data aggregation, content generation, and distribution. Established news agencies such as the Associated Press and BBC use template-driven NLG to automate routine financial and sports reports (Carlson, 2015; Nguyen and Hekman, 2022). Other organizations leverage machine learning algorithms for tasks like topic clustering, headline optimization, and personalized recommendations. Chapter 4s simulation pipeline confirmed these roles by showing how scraped data can be passed through NLP modules and ultimately rendered into draft articles using an NLG engine (Smith and Jones, 2024).
2. **RQ2: What roles do technologies such as NLP, NLG, and ML play in these implementations?** Our analysis shows that different AI subfields serve distinct functions in journalism. NLP techniques (named-entity recognition, sentiment analysis, etc.) are used for information extraction and content tagging, enabling journalists to process large datasets efficiently. Machine learning algorithms power predictive tasks such as trend analysis, audience segmentation, and automated fact-checking. Natural Language Generation plays the role of drafting readable copy: for example, rule-based NLG systems convert structured data (e.g., sports statistics or financial results) into coherent news narratives (Grimme, 2024; Lermann Henestrosa et al., 2022). In Chapter 4s prototype, NLP extracted key events from news feeds, while NLG generated an initial article draft that was then reviewed and refined by a human editor.
3. **RQ3: What are the key benefits and risks associated with AI in journalism, particularly in relation to editorial control and transparency?** The thesis identified numerous benefits and risks (Chapters 3 and 5). Benefits include improved efficiency, consistency, and the ability to cover data-intensive topics rapidly. For instance, automated systems can generate weather forecasts or summarize sports results in seconds, freeing human reporters for in-depth investigative or analytical work. However, we also found significant risks. These include potential errors or biases

in automated reports, loss of nuanced storytelling, and diminished editorial oversight (Steiner, 2019). Transparency issues are prominent: black-box algorithms can obscure how a story was constructed, raising accountability concerns. Chapter 5s discussion highlighted that maintaining editorial control and clear human-review processes is crucial to mitigate these risks.

4. **RQ4: How can generative AI enhance or compromise journalistic standards and integrity?**

This question was addressed by examining advances in generative AI (Chapter 3) and testing a prototype system (Chapter 4). We observed that generative models can assist journalists by proposing leads, suggesting story angles, and translating raw data into narrative prose, thus accelerating content creation. Our prototype incorporated a large language model and showed how it could propose sections of a news article from structured inputs. However, generative AI also carries risks: models may produce hallucinations (false or fabricated information) if not properly constrained (Marcus, 2020; Montoro Montarroso et al., 2023). These hallucinations could compromise integrity if not caught by editors. Our findings suggest that generative AI can enhance productivity when used under human supervision, but it should not replace critical editorial judgment.

5. **RQ5: What design principles can guide the responsible integration of AI into newsroom infrastructures?**

In Chapters 4 and 5, we proposed and evaluated design principles based on ethical and practical considerations. Key principles include transparency (e.g., clearly indicating AI-generated content), editorial control (ensuring humans retain final decision-making authority), data privacy compliance (adhering to GDPR and related laws), and system auditability (logging algorithmic decisions for review) (Rahwan et al., 2019). Our prototype pipeline was designed with modular, open components to allow auditing and updates. We also emphasized the importance of diverse training data and algorithmic accountability to avoid embedding biases (Jobin et al., 2019; Raji and Buolamwini, 2019). These design guidelines provide a framework for responsibly integrating AI into newsroom workflows.

These conclusions underscore recurring themes of bias, editorial control, and transparency, which we examine further in the context of ethical and regulatory considerations below.

6.2 Ethical and Regulatory Considerations

6.2.1 Bias, Transparency, and Trust

One critical concern is the potential for algorithmic bias and reduced transparency. AI systems often learn from large datasets that may contain social biases or historical imbalances (Barocas and Selbst, 2016; O’Neil, 2016). In automated news, biased training data can lead to skewed coverage or misrepresentation of topics. For example, if an NLP model is trained primarily on sources from Western media, it may under-emphasize non-Western perspectives (Diakopoulos, 2025). Such biases threaten journalistic impartiality and diversity, and they can exacerbate existing inequities. Bias can also emerge from the selection of data sources and the design of algorithms. For instance, if an AI system is trained predominantly on wire-service content, it may prioritize official narratives while overlooking local or dissenting voices. Gender or racial biases in text corpora can influence the tone and framing of news stories (Rao and Taboada, 2021). Some newsrooms use fairness audits to detect these biases (Raza et

al., 2022). However, completely eliminating bias is difficult, and thus human oversight remains essential to ensure accuracy. Algorithmic opacity is another major issue. Many AI models, especially complex deep learning systems, operate as black boxes that are difficult for humans to interpret (Burrell, 2016). In practice, this can leave journalists and editors uncertain about how a news story was generated or why a particular article was prioritized. The literature argues for incorporating explainable AI (XAI) and audit logs into journalistic tools so that outputs can be traced and validated. Transparency in the AI pipeline is essential for maintaining trust: journalists need to verify the facts and reasoning behind AI-assisted content. Generative AI introduces the specific risk of hallucinations, where models produce plausible but false or fabricated content. In newsrooms, a hallucinated quote or statistic could propagate misinformation if not caught by editors. Chapter 4s prototype highlighted this danger: when fed ambiguous inputs, the GenAI component sometimes invented fictitious context or figures. Ensuring that AI outputs are factually grounded requires strong guardrails. Current debates stress that explainable AI techniques and integrated fact-checkers should be used to flag potential hallucinations (J. Li et al., 2016). These technical issues directly impact editorial standards and public trust. Journalism relies on accountability and rigorous fact-checking; opaque AI processes can erode these foundations (Radcliffe, 2025). For readers to trust AI-assisted news, media organizations must clearly disclose the use of AI and accept responsibility. For example, survey research indicates that labeling content as AI-generated can significantly affect audience trust, highlighting the need for transparency (Guo, 2022). Researchers have suggested measures such as auditing algorithmic decision processes and enabling post-publication reviews. Ultimately, trust will depend on keeping humans in the loop and ensuring that AI serves rather than supplants human editorial judgment.

6.2.2 The European Legislative Landscape (AI Act, GDPR, etc.)

In the European context, data protection laws have significant implications for AI-driven journalism. The General Data Protection Regulation (GDPR) governs the collection and processing of personal data, which affects how news organizations can use user analytics and personalized algorithms (Kaminski, 2019; Wachter et al., 2017). Journalistic use of digital traces (e.g., social media comments, user browsing histories) is legally constrained: automated content recommendation systems must either obtain user consent or rely on legitimate interest, and they must comply with GDPR transparency requirements (Regulation (EU) 2016/679, Recitals 7172). Article 22 of the GDPR introduces a right against fully automated decisions; this could apply if an AI fully determines editorial content or audience targeting without human intervention. As a result, newsrooms must ensure sufficient human oversight and may need to provide explanations for algorithmic choices, aligning with data subject rights. Notably, the GDPR includes a limited press exemption (Article 85) intended to balance data rights with freedom of expression, but its scope is ambiguous. The right to be forgotten and data minimization requirements may conflict with journalistic archives and data journalism projects. These tensions highlight the need to interpret regulations carefully to avoid hindering press freedom. More recently, the proposed EU AI Act introduces a framework for high-risk AI applications. While journalism is not explicitly listed, certain uses such as automated editing or news curation could fall under high-risk if they significantly affect public discourse (Cabrera, 2024). The draft Act requires providers of AI systems to conduct risk assessments, maintain documentation, and ensure transparency. For example, systems that moderate content or profile readers may be classified as high-risk, imposing strict obligations. This approach

aims to protect against harmful disinformation, but it also creates compliance challenges. News organizations must balance the need for innovation with the workload of compliance procedures, which could disadvantage smaller outlets (European Commission, 2025). Another tension arises between editorial autonomy and regulatory oversight. Journalists traditionally make independent decisions about news coverage, but algorithmic tools inevitably influence these decisions. Regulations like the AI Act emphasize accountability, which may constrain editorial discretion. For instance, a personalized news recommender touches on profiling: under the GDPR and forthcoming ePrivacy rules, readers must be informed and allowed to opt out of profiling (Kaminski, 2019). This could constrain targeted content strategies. On the other hand, compliance requirements such as mandated audit trails for AI decisions could improve editorial transparency and public trust. Overall, the regulatory landscape presents a balance: fostering transparency and protecting rights on one side, while potentially slowing the pace of AI adoption in newsrooms on the other.

6.3 Final Thoughts and Conclusion

The societal role of AI in journalism is multifaceted and evolving. On one hand, AI holds promise to enhance news coverage and democratize information access by enabling outlets to serve readers at scale. For instance, small local newsrooms could use AI to auto-generate community updates, addressing the news desert problem (Quéré, 2022). At the same time, there is caution that unchecked automation could accelerate misinformation and disrupt journalistic jobs. This thesis argues that a balanced, hybrid model is the most viable path forward: AI should augment human journalists rather than replace them. In the envisioned hybrid future, human expertise and AI capabilities complement each other. Journalists bring creativity, contextual awareness, and ethical judgment; AI contributes data-processing speed and consistency. Already, media companies are experimenting with centaur models where AI generates basic drafts and humans edit and enrich the content. Our findings support this approach: news workflows that maintain human oversight at key points tend to preserve quality while benefiting from AI efficiency. Ethical design is crucial: systems should be built with fairness, accountability, and transparency in mind. Newsrooms may need new roles such as algorithmic auditors or data editors to bridge the technical and editorial domains. Continual training and adaptive processes will help journalists stay informed about AI capabilities and limitations. Several future research directions arise from this work. These include:

- Developing reinforcement learning frameworks that incorporate continuous editorial feedback to refine AI outputs).
- Designing robust fact-checking and source-verification mechanisms within AI news pipelines to detect and prevent misinformation.
- Conducting cross-jurisdiction policy studies to compare how different regulations (such as the EU AI Act versus US guidelines) affect news innovation and rights protection.
- Investigating audience trust and engagement with AI-assisted news, including experiments on labeling and transparency effects.
- Evaluating the long-term impacts of AI adoption on newsroom labor structures and journalist roles).

- Exploring federated learning and privacy-preserving techniques to support collaborative journalism models while protecting user privacy.

In conclusion, this thesis contributes to understanding how AI can be responsibly integrated into journalism. The analysis suggests that while AI can improve efficiency and open new possibilities, it also raises ethical and practical challenges that must be addressed through thoughtful design and regulation. A hybrid humanAI approach, underpinned by transparency and accountability, offers the most promise. Journalism's core mission to inform the public with accurate, trustworthy news remains central. As technology evolves, ongoing dialogue among journalists, technologists, and policymakers will be essential to ensure AI serves the public interest. The frameworks and recommendations provided here lay a foundation for such future work.

Bibliography

- Altay, S., & Gilardi, F. (2024). People are skeptical of headlines labeled as ai-generated, even if true or human-made, because they assume full ai automation. *PNAS Nexus*, 3(10), pgae403. <https://doi.org/10.1093/pnasnexus/pgae403>
- Appelman, A., & Bien-Aimé, S. (2024). When experts attack: Readers trust in ai-generated news. *Journalism Studies*, 25(3), 245–262.
- Bandari, R., Asur, S., & Huberman, B. A. (2012). The pulse of news in social media: Forecasting popularity [Accessed: 2025-04-11]. *Proceedings of the Sixth International AAAI Conference on Weblogs and Social Media (ICWSM)*, 26–33. <https://ojs.aaai.org/index.php/ICWSM/article/view/14261>
- Barca, A. (2022). *Il futuro dell'informazione: L'utilizzo dell'intelligenza artificiale nel giornalismo* [Master's thesis, Luiss Guido Carli]. https://tesi.luiss.it/35951/1/646082_BARCA_ALESSIO.pdf
- Barocas, S., & Selbst, A. D. (2016). Big datas disparate impact. *California Law Review*, 104, 671–732.
- Boudet, J., & Vollhardt, K. (2023). What is personalization? <https://www.mckinsey.com/featured-insights/mckinsey-explainers/what-is-personalization>
- Bradshaw, P., & Rohumaa, L. (2011). *Data journalism: Mapping the future*. Sveriges Radio Digital Knowledge Journal.
- Brown, S., Smith, J., & Williams, E. (2020). Automated journalism: The effects of ai authorship and evaluative feedback on news credibility. *Computers in Human Behavior*, 106, 106267. <https://doi.org/10.1016/j.chb.2020.106267>
- Burrell, J. (2016). How the machine thinks: Understanding opacity in machine learning algorithms. *Big Data & Society*, 3(1), 2053951715622512. <https://doi.org/10.1177/2053951715622512>
- Cabrera, L. L. (2024, May). Eu ai act brief pt. 3, freedom of expression [Center for Democracy and Technology]. <https://cdt.org/insights/eu-ai-act-brief-pt-3-freedom-of-expression/>
- Carlson, M. (2015). The robotic reporter: Automated journalism and the redefinition of news. *Digital Journalism*, 3(3), 416–431. <https://doi.org/10.1080/21670811.2014.976412>
- Center, P. R. (2024). News consumption and platforms fact sheet. *Pew Research Center*.
- Choi, J. R. (2018). Newspaper journalists' attitudes towards robot journalism [Accessed: 2025-04-23]. *Computers in Human Behavior*, 85, 72–82. <https://doi.org/10.1016/j.chb.2018.03.037>
- Dalgali, A., & Crowston, K. (2020). Algorithmic journalism and its impacts on work [Accessed: 2025-04-17]. *Proceedings of the Computation + Journalism Symposium*. https://bpb-us-e1.wpmucdn.com/sites.northeastern.edu/dist/d/53/files/2020/02/CJ_2020_paper_26.pdf
- Danzon-Chambaud, S. (2021). A systematic review of automated journalism scholarship: Guidelines and suggestions for future research [Version 1; peer review: 2 approved]. *Open Research Europe*, 1(4). <https://doi.org/10.12688/openreseurope.13096.1>

- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. <https://arxiv.org/abs/1810.04805>
- Diab, K. (2023). *What future for journalism in the age of ai?* [Accessed: 2025-04-17]. <https://www.aljazeera.com/opinions/2023/7/19/what-future-for-journalism-in-the-age-of-ai>
- Diakopoulos, N. (2019). *Automating the news: How algorithms are rewriting the media*. Harvard University Press.
- Diakopoulos, N. (2025). Prospective algorithmic accountability and the role of the news media [Available at SSRN: <https://ssrn.com/abstract=5022167>]. In M. Noorman & M. Verdicchio (Eds.), *Computer ethics across disciplines: Algorithmic accountability and ai through deborah g. johnsons lens*. Springer. <https://doi.org/10.2139/ssrn.5022167>
- Dörr, K. (2016). Mapping the field of algorithmic journalism [Accessed: 2025-03-02]. *Digital Journalism*, 4(6), 700–722. <https://doi.org/10.1080/21670811.2016.1188460>
- Dörr, K. N. (2015). Mapping the field of algorithmic journalism. *Digital Journalism*, 4(6), 700–722. <https://doi.org/10.1080/21670811.2015.1096748>
- Enevoldsen, K., Chung, I., Kerboua, I., Kardos, M., Mathur, A., Stap, D., Gala, J., Siblini, W., Krzemiski, D., Winata, G. I., Sturua, S., Utpala, S., Ciancone, M., Schaeffer, M., Sequeira, G., Misra, D., Dhakal, S., Rystrom, J., Solomatin, R., ... Muennighoff, N. (2025). Mmteb: Massive multilingual text embedding benchmark. <https://arxiv.org/abs/2502.13595>
- European Commission. (2025). European approach to artificial intelligence [Accessed: 2025-05-15]. <https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>
- Fanta, A. (2017). *Putting europe's robots on the map: Automated journalism in news agencies* (tech. rep.). Reuters Institute for the Study of Journalism. <https://reutersinstitute.politics.ox.ac.uk/sites/default/files/2017-09/Fanta%2C%20Putting%20Europe%E2%80%99s%20Robots%20on%20the%20Map.pdf>
- Fearn, N. (2025, April). *How to prevent ai from exacerbating diversity and inclusion in the journalism industry* [Accessed: 2025-04-04]. <https://www.media-diversity.org/how-to-prevent-ai-from-exacerbating-diversity-and-inclusion-in-the-journalism-industry/>
- Graefe, A. (2016). *Guide to automated journalism* (tech. rep.). Tow Center for Digital Journalism, Columbia University. <https://academiccommons.columbia.edu/doi/10.7916/D80G3XDJ>
- Grimme, M. (2024). *Ai in media organisations: Factors influencing the integration of ai in the newsroom* [Doctoral dissertation, University of Hohenheim]. <https://hohpublica.uni-hohenheim.de/bitstreams/91188871-64a3-4f2a-b6ad-c5cf6c017e88/download>
- Guo, Y. (2022). Audience trust in ai-generated news: An experimental study. *Journalism and Media*.
- Hasan, K., Kulkarni, R., Bijale, M., Naik, N., Bhat, P., & Kulkarni, M. (2023). Impact of ai integration on journalists mental health: A quantitative study [Accessed: 2025-04-12]. *International Journal of Communication and Health (IJCH)*, 1(2). <https://doi.org/10.5547/10.00.XXXXX>
- Jiang, Y. (2020). Sst-2 sentiment analysis [Accessed May 10, 2025].
- Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of ai ethics guidelines. *Nature Machine Intelligence*, 1(9), 389–399. <https://doi.org/10.1038/s42256-019-0088-2>
- Johnson, M., & Black, R. (2020). Who owns the output? ethics of responsibility in algorithmic journalism. *AI and Ethics*.

- Kaminski, M. E. (2019). The right to explanation, explained. *Berkeley Technology Law Journal*, 34(1), 189–218.
- Kumar, S. D. M., Athira, A. B., & Chacko, A. M. (2022). Towards smart fake news detection through explainable ai [Accessed: 2025-03-28]. *arXiv preprint arXiv:2207.11490*. <https://doi.org/10.48550/arXiv.2207.11490>
- Leppänen, L., Munezero, M., Granroth-Wilding, M., & Toivonen, H. (2017, September). Data-driven news generation for automated journalism. In J. M. Alonso, A. Bugarín, & E. Reiter (Eds.), *Proceedings of the 10th international conference on natural language generation* (pp. 188–197). Association for Computational Linguistics. <https://doi.org/10.18653/v1/W17-3528>
- Lermann Henestrosa, A., Greving, H., & Kimmerle, J. (2022). Automated journalism: The effects of ai authorship and evaluative information on the perception of a science journalism article. *Computers in Human Behavior*, 138. <https://doi.org/10.1016/j.chb.2022.107445>
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., & Zettlemoyer, L. (2020). Bart: Denoising sequence-to-sequence pre-training for natural language generation [Accessed: 2025-04-27]. *Proceedings of ACL*, 7871–7880. <https://aclanthology.org/2020.acl-main.703/>
- Li, J., Chen, X., Hovy, E., & Jurafsky, D. (2016). Visualizing and understanding neural models in NLP, 681–691. <https://doi.org/10.18653/v1/N16-1082>
- Li, M., & Wang, L. (2019). A survey on personalized news recommendation technology. *IEEE Access*, 7, 145861–145879. <https://doi.org/10.1109/ACCESS.2019.2944927>
- Linden, C.-G., Kalogeropoulos, A., & Nielsen, R. K. (2021). *New powers, new responsibilities: A global survey of journalism and artificial intelligence* (tech. rep.) (Accessed: 2025-05-02). Reuters Institute for the Study of Journalism. <https://reutersinstitute.politics.ox.ac.uk/new-powers-new-responsibilities-global-survey-journalism-and-artificial-intelligence>
- Marconi, F. (2019). *Newsmakers: Artificial intelligence and the future of journalism*. Columbia University Press. <https://cup.columbia.edu/book/newsmakers/9780231191364>
- Marcus, G. (2020). *Rebooting ai: Building Artificial Intelligence we can trust*. Pantheon.
- McFarland, M. (2015, January). *Associated press looks to expand its automated stories program following successful launch*. <https://www.washingtonpost.com/news/innovations/wp/2015/01/29/associated-press-looks-to-expand-its-automated-stories-program-following-successful-launch/>
- Montoro Montarroso, A., Cantón, J., Rosso, P., Chulvi, B., Panizo Lledot, Á., Huertas-Tato, J., Figueras, B., Rementeria, M., & Gómez-Romero, J. (2023). Fighting disinformation with artificial intelligence: Fundamentals, advances and challenges. *El Profesional de la información*, 32. <https://doi.org/10.3145/epi.2023.may.22>
- Narayanan, A. (2018, October). How to recognize ai snake oil [Presentation at MIT’s Program in Science, Technology, and Society (STS)]. <https://www.cs.princeton.edu/~arvindn/talks/MIT-STS-AI-snakeoil.pdf>
- NewsGuard. (n.d.). Tracking ai-enabled misinformation: Over 1200 ‘unreliable ai-generated news’ websites (and counting), plus the top false narratives generated by artificial intelligence tools [Accessed: 2025-04-07]. <https://www.newsguardtech.com/special-reports/ai-tracking-center/>
- Nguyen, D., & Hekman, E. (2022). The news framing of artificial intelligence: A critical exploration of how media discourses make sense of automation. *AI & SOCIETY*, 39. <https://doi.org/10.1007/s00146-022-01511-1>

- ONeil, C. (2016). *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown Publishing.
- Pariser, E. (2011). *The filter bubble: What the internet is hiding from you*. Penguin Press.
- Posetti, J. (2018, November). *Time to step away from the 'bright, shiny things'? towards a sustainable model of journalism innovation in an era of perpetual change* (tech. rep.) (Accessed: 2025-04-09). Reuters Institute for the Study of Journalism. https://reutersinstitute.politics.ox.ac.uk/sites/default/files/2018-11/Posetti_Towards_a_Sustainable_model_of_Journalism_FINAL.pdf
- Press, A. (2016). *Ap expands minor league baseball coverage* [Accessed: 2025-03-18]. <https://www.ap.org/media-center/press-releases/2016/ap-expands-minor-league-baseball-coverage>
- Press, A. (2024). Artificial intelligence initiatives. *AP.org*.
- Puccetti, G., Rogers, A., Alzetta, C., Dell'Orletta, F., & Esuli, A. (2024). Ai "news" content farms are easy to make and hard to detect: A case study in italian [In proceedings of ACL 2024]. *arXiv preprint arXiv:2406.12128*. <https://arxiv.org/abs/2406.12128>
- Quéré, M. A. L. (2022, September). Rethinking ai and local news [Digital Life Initiative, Cornell Tech]. <https://dli.tech.cornell.edu/post/rethinking-ai-and-local-news>
- Radcliffe, D. (2025). *Journalism in the ai era: Opportunities and challenges in the global south and emerging economies* (tech. rep.). Thomson Reuters Foundation.
- Rahwan, I., Cebrian, M., Obradovich, N., Bongard, J., Bonnefon, J.-F., Breazeal, C., Crandall, J., Christakis, N., Couzin, I., Jackson, M., Jennings, N., Kamar, E., Kloumann, I., Larochelle, H., Lazer, D., McElreath, R., Mislove, A., Parkes, D., Pentland, A., & Wellman, M. (2019). Machine behaviour. *Nature*, 568, 477–486. <https://doi.org/10.1038/s41586-019-1138-y>
- Raji, I., & Buolamwini, J. (2019). Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products, 429–435. <https://doi.org/10.1145/3306618.3314244>
- Rao, P., & Taboada, M. (2021). Gender bias in the news: A scalable topic modelling and visualization framework. *Frontiers in Artificial Intelligence*, 4, 664737. <https://doi.org/10.3389/frai.2021.664737>
- Raza, S., Reji, D., & Ding, C. (2022). Dbias: Detecting biases and ensuring fairness in news articles. *International Journal of Data Science and Analytics*, 17, 39–59. <https://doi.org/10.1007/s41060-022-00359-4>
- Sawers, P. (2014). *How the associated press is using automation to produce thousands of earnings reports*. <https://venturebeat.com/2014/07/11/how-the-associated-press-is-using-automation-to-produce-thousands-of-earnings-reports/>
- Smith, J., & Jones, E. (2024). From template to transformer: Evolving models in automated journalism. *AI & Society*.
- Steiner, L. (2019). *Journalism in crisis: Risks and responsibilities in the digital age*. Polity Press.
- Su, B., & Xu, X. (2019). Understanding the adoption of artificial intelligence in journalism: A study of chinese journalists. *SAGE Open*, 9(3), 215824401986909. <https://doi.org/10.1177/215824401986909>
- The Washington Post. (2020, October). *The washington post to debut ai-powered audio updates for 2020 election results*. <https://www.washingtonpost.com/pr/2020/10/13/washington-post-debut-ai-powered-audio-updates-2020-election-results/>

- Upadhyay, A., Bijale, M., & Hasan, K. (2024). Impact of ai integration on journalists' mental health: A quantitative study [Epub ahead of print]. *Annals of Neurosciences*. <https://doi.org/10.1177/09727531241278909>
- van Dalen and, A. (2012). The algorithms behind the headlines. *Journalism Practice*, 6(5-6), 648–658. <https://doi.org/10.1080/17512786.2012.667268>
- Vega, N. (2021, August). *The washington posts robot reporter is rewriting headlines to suit your interests*. <https://www.theverge.com/2021/8/5/22610188/washington-post-heliograf-ai-headlines-newsroom>
- Wachter, S., Mittelstadt, B., & Floridi, L. (2017). Why a right to explanation of automated decision-making does not exist in the general data protection regulation. *International Data Privacy Law*, 7(2), 76–99.
- Walker, M. (2021). U.s. newsroom employment has fallen 26% since 2008. *Pew Research Center*.
- WAN-IFRA. (2023). New genai survey [Accessed: 2025-03-02]. <https://wan-ifra.org/2023/05/new-genai-survey/>
- Weedon, G. (2021). Signal ai c-suite whitepaper.
- Wu, C., Wu, F., Huang, Y., & Xie, X. (2021). Personalized news recommendation: Methods and challenges. *arXiv preprint arXiv:2106.08934*. <https://arxiv.org/abs/2106.08934>
- Xu, Z., & Lan, X. (2020). A scientometric review of automated journalism: Analysis and visualization. *Journal of Physics: Conference Series*, 1684(1), 012127. <https://doi.org/10.1088/1742-6596/1684/1/012127>
- Zhang, J., Dong, B., & Yu, P. S. (2019). Fakedetector: Effective fake news detection with deep diffusive neural network. <https://arxiv.org/abs/1805.08751>